

Hybster

A Highly Parallelizable Protocol for Hybrid Fault-Tolerant Service Replication

Johannes Behl

TU Braunschweig

`behl@ibr.cs.tu-bs.de`

Tobias Distler

FAU Erlangen-Nürnberg

`distler@cs.fau.de`

Rüdiger Kapitza

TU Braunschweig

`rrkapitz@ibr.cs.tu-bs.de`

April 21st, 2017
v0.1

Contents

1	Introduction	1
2	Notation	1
2.1	Modules and Operations	1
2.2	Module Interaction	3
2.3	Additional Expressions and Types	6
3	System Model	8
3.1	Processes	8
3.2	Time	10
3.3	Communication	11
3.4	Cryptography	11
3.5	Faults	13
4	Specification of Hybster	15
4.1	System Properties	15
4.2	Auxiliary Modules	17
4.2.1	Adaptive Timeouts	18
4.2.2	Connections	19
4.3	TrInX	21
4.4	Processes	22
4.4.1	Clients	22
4.4.2	Replicas	22
4.5	Invocation Protocol	24
4.5.1	Client Side	24
4.5.2	Replica Side	25
4.6	Ordering Protocol	27
4.7	Execution	31
4.8	Checkpointing Protocol	32
4.9	View-Change Protocol	35
4.10	State Transfer Protocol	41

1 Introduction

Hybster is a replication protocol for providing fault-tolerant services based on the state-machine approach [13]. It is designed for partially synchronous systems [9]. It does not rely on timing assumptions to ensure safety but is only guaranteed to make progress in phases where unknown upper bounds for message delivery and message processing exist. Moreover, using a hybrid fault model in which some components are trusted to only fail by crashing even when other components behave arbitrarily faulty [12, 14], Hybster requires only $2f + 1$ service replicas to tolerate up to f faulty ones.

While the original paper that introduces Hybster [3] mainly presents the rationale behind the basic protocol and its parallelized variant, this technical report provides supplemental material such as a comprehensive system model description and formal specification of the protocol. The present document is the first version of the report that is planned to be updated once additional parts are ready. All published versions will be accessible at [2].

2 Notation

The formal specification of Hybster is essentially based on I/O automata [11] as employed by Castro and Liskov to specify the originator of Byzantine fault-tolerant service replication, namely PBFT [5, 6, 8]. However, the model used here differs from I/O automata mainly in two aspects. Firstly, it supports the direct interaction of modules where module instances explicitly trigger operations on other module instances. This makes dependencies between modules more apparent than the indirect interaction scheme based on operation signatures as used by the original I/O automata. Moreover, it allows more flexible interaction patterns between modules, especially helpful for describing the internal structure of complex processes. Secondly, to improve readability, the notation used in this document is inspired by pseudo-code notations like [4].

2.1 Modules and Operations

The specification of Hybster is given as a set of *modules* encapsulating *state* and defining *operations* on that state. Compared to an I/O automaton, a module is a more abstract concept. A module is not only a concurrent entity but allows for both asynchronous and synchronous operations. Furthermore, modules can be instantiated, there can be multiple identifiable *instances* of one and the same module. Instances share the behavior as described by the operations of the module but possess independent memory for state variables. Thus, modules and module instances have strong similarities to classes and objects of the object-oriented programming paradigm [1].

Tasks and Methods. Asynchronous operations are called *tasks* and are triggered according to a specified precondition. If the precondition of a task is satisfied for a module instance, the task of this instance is *enabled* and all enabled tasks are assumed to be executed asynchronously at some time, in a nondeterministic order, although only one task at a time. Tasks are classified as input, output, or internal operation. While tasks are not allowed to yield a direct return value, output tasks may specify a list of output parameters that serve as inputs for other tasks. Likewise, input tasks define a list of expected input parameters. Internal tasks can possess a parameter list as well. These parameters are neither inputs nor outputs and their only purpose is to make the function, the logical result of the internal tasks explicit.

Synchronous operations, called *methods*, are directly executed when invoked and may return a direct outcome. Methods are also classified as input, output, and internal operation but can be additionally declared as public. A public method is permitted to alter the state of the invoked module instance and yield a direct result; a public method can be regarded as both input and output. Besides tasks and methods, *global functions* can be specified. Such operations are not bound to instances of modules thus they cannot access any instance state.

Figure 1: Example for module definitions.

```

1 module NumberProducerp
2   upon init() do
3     outp ∈ (ℕ)* := []

5   upon internal task produce(val ∈ ℕ)
6     with
7       |outp| < 10
8       val = random(0, 100)
9     do
10      outp.append(val)

12 module NumberConsumerc
13   upon init() do
14     sumc ∈ ℕ := 0

16   upon internal call consume(val ∈ ℕ) do
17     sumc := sumc + val

```

Figure 1 gives an example for the definition of two modules. The module *NumberProducer* generates random numbers within an internal task *produce* and *NumberConsumer* summarizes numbers within an internal method *consume*. Both modules are not yet connected. How they can interact is described in the next section. The freely chosen index *p* in “**module** NumberProducer_{*p*}” identifies an instance of the *NumberProducer* in the context of the module definition. The special method *init* is used to initialize a newly created instance of a module. “*out_p* ∈ (ℕ)* := []” in Line 3 defines a single variable *out* for a new instance of the module *NumberProducer* denoted with *p*. *out_p* is an instance (∈) of a list ((...)*) of natural numbers (ℕ) that is initially assigned (:=) to an empty list ([]). The definition of the *produce* task starts in Line 5. It is not allowed to have a direct return value. However, the parameter *val* is a free variable that is eventually bound to the value *produce* generates. *val* cannot be used directly. At this point, its only purpose is to make the function of *produce*, its logical outcome explicit. The **with** keyword starts the precondition section of the task definition. *produce* is only enabled when the list *out_p* contains fewer than 10 elements. Since *val* is a free variable, binding *val* to some randomly chosen natural number by “*val* = *random*(0, 100)” is always satisfied. In this example, *random* is assumed to be a global function defined within some other module. When a task is enabled, eventually all imperative statements in the section introduced by the keyword **do** will be executed in the sequence as specified. *produce* comprises only a single statement. “*out_p*.*append*(*val*)” adds the value bound to the variable *val* to the end of the list referenced by *out_p*. In sum, *produce* appends randomly chosen numbers to the list *out_p* as long as *out_p* currently contains fewer than 10 elements. The definition of *NumberConsumer* starting in Line 12 uses *c* to refer to an instance of that module. Instances of *NumberConsumer* possess one member variable *sum_c* that is a natural number initialized with 0. Further they possess one internal method *consume* that takes a value as argument and adds the value to *sum_c*.

Operation Definition and Declaration. The complete syntax for the definition of operations can be found in Figure 2. The keyword **upon** begins the definition of a task or a method. It is followed by an access modifier and an optional declaration of the operation type (**task** for tasks and **call** for methods). Operations may define a list of parameters. In the case of output and internal tasks, the parameters are output, in all other cases input parameters. Methods

Figure 2: Syntax for operation definitions.

(a) Define tasks.	(b) Define methods and global functions.
<pre> 1 upon <internal input output> [task] taskName($arg_0 \in T_0, \dots, arg_{n-1} \in T_{n-1}$) 2 asserts 3 // assertions 4 with 5 // preconditions 6 assert $O()$ // included assertion 7 [if] $P()$ 8 do 9 // statements 10 yields 11 // postconditions </pre>	<pre> 1 upon <internal input output public global> [call] methodName($arg_0 \in T_0, \dots$) $\rightarrow ret \in R$ 2 asserts 3 // assertions 4 do 5 // statements 6 return ret 7 yields 8 // postconditions </pre>

may additionally define a direct return parameter, although this is optional. Access modifier, type declaration, name, and the list of all parameters form the *signature* of an operation. It follows a number of sections that constitute the body of the operation definition. *Assertions* are predicates that are expected to be true, they are assumptions that shall be made explicit. Assertions can be declared within the **asserts** section or anywhere else following an **assert**. A task must specify a set of *preconditions* that determine when the task is enabled and thus eligible for execution. Preconditions are given in the **with** section. Single preconditions can be prefaced by an optional **if** and if multiple conditions are specified, they are implicitly combined by a logical AND operation (\wedge). If an enabled task is executed or if a method is invoked, the sequence of *statements* within the **do** section is evaluated. Unless the operation is a global function, statements are allowed to modify the state of the invoked module instance. If a **return** statement is reached, the evaluation of the operation ends. For methods with return parameter, the **return** keyword has to be followed by an expression that specifies the return value. In conjunction with or as replacement for statements that explicitly alter state or define the values of output and return parameters, a set of *postconditions* may be declared that have to be satisfied after the evaluation of an operation. Postconditions are located in the **yields** section.

Figure 3: Syntax for operation declarations.

(a) Declare tasks.	(b) Declare methods and global functions.
<pre> 1 declare <internal input output> task taskName($arg_0 \in T_0, \dots, arg_{n-1} \in T_{n-1}$) </pre>	<pre> 1 declare <internal input output public global> call methodName($arg_0 \in T_0, \dots$) $\rightarrow ret \in R$ </pre>

Operation definitions are both a declaration and a direct implementation. If just a type or a part of an interface is to be given, the syntax shown in Figure 3 can be used to merely declare operations. The implementation has to be provided by a different module.

2.2 Module Interaction

Modules are not only used to describe a single algorithm, they also serve as interacting components in the specification of more complex systems. To describe various patterns of how modules can interact, the modules *NumberProducer* and *NumberConsumer* defined in Figure 1 of the previous section are extended in different ways. In contrast to the *indirect interaction* scheme of I/O automata in which output operations of one automaton are linked to input operations

Figure 4: Direct unidirectional module interaction.

(a) Source-linked.

```

1 module LinkingNumberSourcep
2   extends NumberProducer

4   upon init() do
5     snkp := NumberSink()

7   upon internal task transmit()
8     with
9       |outp| > 0
10    do
11      output snkp.receive(outp.dequeue())

13 module NumberSinkc
14   extends NumberConsumer

16   upon input call receive(val ∈ ℕ) do
17     consume(val)

```

(b) Sink-linked.

```

1 module NumberSourcep
2   extends NumberProducer

4   upon output task transmit(val ∈ ℕ)
5     with
6       |outp| > 0
7     do
8       val := outp.dequeue()

10 module LinkingNumberSinkc
11   extends NumberConsumer

13   upon init() do
14     srcc := NumberSource()

16   upon internal task receive(val ∈ ℕ)
17     with input srcc.transmit(val) do
18       consume(val)

```

of other automata on the basis of the operations' signature [11], in the model employed here, instances of modules invoke operations of other instances *directly*.

Unidirectional Interaction. For example, Figure 4 shows two forms of a *direct unidirectional interaction*. Generally, in an unidirectional interaction, the output of one module, the source, is taken as input for another, the sink. Figure 4a presents the *source-linked* variant of this interaction scheme. In this example, two modules are defined, *LinkingNumberSource* and *NumberSink*. The module *LinkingNumberSource* is derived from the module *NumberProducer* and the current instance of this module is referred to as *p*. *LinkingNumberSource* inherits from *NumberProducer* the instance variable *out_p* maintaining the list of generated but not yet consumed random numbers and the internal task *produce*. Additionally, it defines the instance variable *snk_p* initialized with a new instance of the module *NumberSink* and an internal task *transmit*. *transmit* is enabled when *out_p* contains at least one generated number. Upon being executed, *transmit* removes the first value from the list *out_p* (*out_p.dequeue*()) and invokes the input method *receive* of the instance *snk_p* with this value as argument. The method invocation is synchronous, that is, all of its effects take place immediately. While being optional, the preceding keyword **output** shall emphasize this particular interaction between *LinkingNumberSource* and *NumberSink*. The module *NumberSink* is derived from *NumberConsumer* and thus inherits the variable *sum_c* and the internal method *consume*, with *c* denoting the current module instance. The input method *receive* of *NumberSink* simply invokes *consume* of its base module handing over the given value *val* as argument. Taken all together, the task *produce* defined by *NumberProducer* asynchronously generates random natural numbers and places them in the list *out_p* as long as this list contains fewer than 10 values. The task *transmit* asynchronously removes values from this list and delivers them as input to the *NumberSink* instance *snk_p*. The invoked input method *receive* of *snk_p* calls the internal method *consume* that finally adds the given value to the instance variable *sum_c*, with *c* being the internal reference of the module instance also referenced by *snk_p*.

Taken a source-linked interaction, the source is bound to the type of the sink and the source

instance needs a reference to the instance of the sink. In a *sink-linked* interaction scheme, this relation is reversed. The example depicted in Figure 4b resembles the previous one. The module *NumberSource* extends *NumberProducer* and defines the operation *transmit* whereas the module *LinkingNumberSink* extends *NumberConsumer* and defines the operation *receive*. However, in this case, *transmit* is an output task that delivers the next random value from the list out_p by means of the output parameter val . The addressee of this value is not known to the *NumberSource* module. Instead, instances of *LinkingNumberSink* maintain a reference src_c to a *NumberSource* instance and their internal *receive* task is bound to $src_c.transmit$, that is, *receive* is directly called after $src_c.transmit$ has been executed. Analogous to the **output** keyword, the **input** keyword emphasizing the interaction between *NumberSource* and *LinkingNumberSink* may be omitted. While receiving operations in sink-linked interactions are allowed to wait for multiple output events from different instances or even different output operations, they must not specify additional preconditions that would prevent any of their awaited outputs from being executed. In other words, sink-linked interactions must remain input-enabled as required from I/O automata [11].

Bidirectional Interaction. Whereas unidirectional interaction schemes link outputs to inputs, *method invocations* are a form of *bidirectional interaction* in which a callee takes an input and directly returns an output. Figure 5a presents an example. Again, *NumberProducer* and *NumberConsumer* from Figure 1 are extended to create interacting modules. The resulting modules *IncreasingSum* and *Accumulator* resemble the modules *LinkingNumberSource* and *NumberSink* of the source-linked interaction example presented in Figure 4a. Opposed to the input method *receive* of *NumberSink*, the public method *add* of *Accumulator* does not only add the given value to the variable sum_c but also returns the resulting sum as direct outcome of the method invocation. To emphasize the invocation at the caller's side, the keyword **invoke** can be used.

Figure 5: Further module interaction schemes.

(a) Method invocation.	(b) Global function.
<pre> 1 module IncreasingSum_p 2 extends NumberProducer 4 upon init() do 5 $acc_p := Accumulator()$ 6 $sum_p \in \mathbb{N} := 0$ 8 upon internal task transmit() 9 with 10 $out_p > 0$ 11 do 12 $sum_p := \text{invoke } acc_p.add(out_p.dequeue())$ 14 module Accumulator(NumberConsumer)_c 15 upon public call add($val \in \mathbb{N}$) $\rightarrow sum \in \mathbb{N}$ do 16 $consume(val)$ 17 return sum_c </pre>	<pre> 1 module Computer_c 2 upon global call thinkDeep() $\rightarrow \mathbb{N}$ do 3 return 42 5 module Being_b 6 upon init() do 7 $answer_b \in \mathbb{N} \cup \{\emptyset\} := \emptyset$ 9 upon internal task askUltimate() 10 with 11 $answer_b = \emptyset$ 12 do 13 $answer_b := Computer::thinkDeep()$ </pre>

Global Functions. *Global functions* share the syntax with methods but are not bound to any instance. While being specified within some module, they can be called from anywhere. Fig-

ure 5b illustrates the usage. The function *thinkDeep* is defined within the module *Computer* and simply returns 42 as constant. It is called by the module *Being* via “*Computer::thinkDeep()*”. Stating the defining module explicitly is, however, optional.

Inner Modules. Another form of interaction that leads to a very tight coupling between modules is the use of *inner modules*. An instance of an inner module possesses an implicit reference to the module instance that creates it and has access to all instance variables and operations be it internal or external of the referenced outer module instance. Figure 6 gives an example. Instances of the module *OuterNumberSource* create an instance of the inner module *InnerNumberSink*. The inner module instance thereby gets access to the list *out_c* of produced random numbers that *OuterNumberSource* inherits from *NumberProducer*. If *out_c* contains values, the internal task *consumeInner* of *InnerNumberSink* is enabled. When executed, the inner instance removes a value from *out_c* maintained in the outer instance and passes it as argument to the *consume* method *InnerNumberSink* inherits from *NumberConsumer*.

Figure 6: Inner modules.

```

1 module OuterNumberSourcep
2   extends NumberProducer

4   upon init() do
5     snkp := OuterNumberSourcep::InnerNumberSink()

7   module OuterNumberSource::InnerNumberSinkc
8     extends NumberConsumer

10  upon internal task consumeInner(val ∈ ℕ)
11    with
12      |outc| > 0
13    do
14      consume(outc.dequeue())

```

2.3 Additional Expressions and Types

While modules and their various ways of interaction are related to the general architecture of a system, the purpose of operations is to specify concrete algorithms by declaring pre- and postconditions or by defining them step-by-step as sequence of actions. This section gives an overview of expressions and basic types used to form predicates and imperative statements for the definition of operations.

Sets. Sets as unordered collections of objects are denoted by curly brackets. For example, {1, 2, 3} defines the set of the natural numbers 1, 2, and 3. The empty set can be written as {} or ∅. Concrete sets are instances of the basic type *Set*. $X \in \text{Set}$ declares an arbitrary set *X*, and $X := \{1, 2, 3\}$ assigns the set {1, 2, 3} to it. Common set operations are determining the union (∪) or the difference (∖) of two sets. For example, $X := X \setminus \{2\} \cup \{4\}$ results in $X = \{1, 3, 4\}$. The power set of some set *Y* is given by $\mathcal{P}(Y)$. This can be used to specify a typed set: $Z \in \mathcal{P}(\mathbb{N})$ declares *Z* as a set of natural numbers.

Tuples. The ordered counterparts of sets are tuples which are instances of the basic type *Tuple*. Given *n* objects *x*₀ to *x*_{*n*−1}, the corresponding tuple is written as $\langle x_0, \dots, x_{n-1} \rangle$. For

example, $t \in \text{Tuple} := \langle 1, 2, 3 \rangle$ defines t as the tuple comprising the elements 1, 2, and 3 in that order. t could have been also declared as $t \in \mathbb{N} \times \mathbb{N} \times \mathbb{N}$ or $t \in \mathbb{N}^3$, additionally specifying the types of the tuple's elements.

Lists. While tuples are immutable, lists are ordered collections that can be modified. Lists are instances of the basic type *List* that is also denoted by $(\dots)^*$. For example, a list l of natural numbers can be declared with $l \in (\mathbb{N})^*$. Lists are given within square brackets. $l := [1, 2, 3]$ initializes l with the elements 1, 2, and 3. The size of l is $|l| = 3$. The type *List* defines the methods *append*, *pop*, and *dequeue*. *append* adds an object to the end of the list, *pop* returns and removes the last element and *dequeue* the first element. Taken the list l for example, $l.append(4)$ results in $l = [1, 2, 3, 4]$, $x := l.pop()$ in $x = 3 \wedge l = [1, 2]$, and $x := l.dequeue()$ in $x = 1 \wedge l = [2, 3]$.

Booleans. Classical truth values are represented by the basic type $\text{Bool} \stackrel{\text{def}}{=} \{\text{True}, \text{False}\}$, with $\text{True} \stackrel{\text{def}}{=} \emptyset = \emptyset$ and $\text{False} \stackrel{\text{def}}{=} \emptyset \neq \emptyset$. Based on that, a predicate can be defined as a method or function returning a value of type *Bool*: **upon** $[\dots]$ **call** $\text{predicateName}(\dots) \rightarrow \text{Bool}$. An alternative notation for *True* is \top and for *False* \perp .

Ellipses. Predicates may take several input parameters. Sometimes it is sufficient to verify that a predicate holds for a subset of these parameters in conjunction with any binding of the remaining ones. For example, given a ternary predicate P and two particular arguments x and y , the question could be if P is fulfilled for x and y and any additional argument z . To simplify the resulting expression $\exists z P(x, y, z)$, the elliptic form $P(x, y, \cdot)$ can be used. The extended variant for multiple arbitrarily chosen input parameters is written as $P(x, \dots) \stackrel{\text{def}}{=} \exists y \exists z P(x, y, z)$.

None. At some occasions it can be helpful to bind a variable to a special value that resides outside the main domain of that variable, for instance, to denote that the variable is unset and does not contain any real value. As such special value *None* may be used, also denoted by \emptyset . With $x \in \mathbb{N} \cup \{\emptyset\}$, the variable x can reference a natural number or the special value \emptyset .

In addition to basic types, specifications can make use of common control flow structures such as conditionals and loops but also of more specific structures with regard to asynchronous operations.

Sequences of Statements. Sequences of statements are usually given as list, each statement in a single line at the same level of indentation. Figure 7a shows an example in which the values of two variables x and y are swapped using an additional variable tmp . An alternative notation of this example where statements are written in a single line separated by semicolons is given in Figure 7b. If single statements are wrapped, be it for clarity or because they do not fit into the space for a single line, overhanging parts have to be indented by two levels. As Figure 7c illustrates, this can happen within a single line or by splitting the statement such that it spans multiple lines.

Conditionals and Loops. The notation also supports typical control flow statements for conditional branches or repetitive executions. Figure 8a gives an example in which two variable x and y are compared and a return value is chosen depending on the outcome of this comparison. If a condition is satisfied, the associated block indented by one level is executed. Although these

Figure 7: Sequences of statements.

(a) Multiple lines.	(b) Single line.	(c) Long lines.
1 $tmp := x$ 2 $x := y$ 3 $y := tmp$	1 $tmp := x; x := y; y := tmp$	1 $lx :=$ "This is a long expression written in a single line." 2 $ly :=$ "And this is another long expression but written 4 using multiple lines."

blocks contain only a single **return** in this example, blocks can generally span an arbitrary number of statements. An example for a loop that iterates over all elements of a set X and passes them to a procedure $processX$ is depicted in Figure 8b.

Figure 8: Control flow statements.

(a) Conditionals.	(b) Loops.
1 if $x < y$ then 2 return x 3 else if $y < x$ then 4 return y 5 else 6 return 0	1 for all $x \in X$ do 2 $processX(x)$

Synchronous Task Evaluations. As described in Section 2.1, tasks are enabled when their specified precondition is satisfied but they are only executed asynchronously without any guaranteed order. Furthermore, if the precondition of a task t depends on state that is altered by other tasks such that t is enabled and disabled depending on which other tasks are executed, there is no guarantee that the task t is executed at all. To support occasions in which a task is to be executed synchronously if necessary, the **check** keyword is provided. For example, if “**check** $someTask()$ ” is called, the precondition of task $someTask$ is evaluated and, if the precondition is satisfied, $someTask$ is executed like a synchronous operation. Thus, the task is transformed into a conditional method invocation.

3 System Model

The system model Hybster is based on is similar to the one used by PBFT [5, p. 26ff.]. Briefly, in a partially synchronous system with unreliable network, clients invoke operations of a stateful service that is replicated across several servers called replicas in order to ensure that the service does not return incorrect results even if a bounded number of these replicas behaves arbitrarily faulty. As the most notable difference to PBFT, however, Hybster also assumes that processes can be equipped with subsystems that are more reliable than other components such that these subsystems never fail in an arbitrary manner but only by crashing. Assuming such trusted subsystems, Hybster requires fewer replicas than PBFT to tolerate the same number of faults. A more detailed description of Hybster’s system model is given in the following.

3.1 Processes

Hybster is designed for a distributed system in which processes communicate exclusively via message passing over a shared network. Two types of processes are distinguished, *replicas* and

clients. A set of replicas, also called *replica group*, provides a replicated realization of a service, with each replica hosting an instance of an implementation of that service. Clients invoke operations of the service by means of a local invocation handler instance encapsulating the communication with the replicas. The set of replicas in the system $\{r_0, r_1, \dots, r_{n-1}\}$ is denoted by \mathcal{R} and the number of replicas by n , that is, $n \stackrel{\text{def}}{=} |\mathcal{R}|$. Analogously, the set and numbers of clients is given by $\mathcal{C} \stackrel{\text{def}}{=} \{c_0, c_1, \dots, c_{nclients-1}\}$ and $nclients \stackrel{\text{def}}{=} |\mathcal{C}|$.

Listing 1: Specification of a deterministic stateful service implementation.

```

1 module Services

3 upon init() do
4   states ∈  $\mathcal{S} := s_0$ 

6 upon public call invoke( $c \in \mathcal{C}, invno \in {}^i\mathbb{N}, svccmd \in \mathcal{O}$ ) →  $\mathcal{O}'$  do
7    $s', svcret := calculateResult(state_s, c, invno, svccmd)$ 
8   states :=  $s'$ 
9   return svcret

11 declare internal call calculateResult( $s \in \mathcal{S}, c \in \mathcal{C}, invno \in {}^i\mathbb{N}, svccmd \in \mathcal{O}$ ) →  $\mathcal{S} \times \mathcal{O}'$ 

13 upon output call createStateSnapshot() →  $\mathcal{S}$  do
14   return states

16 upon input call installStateSnapshot( $s \in \mathcal{S}$ ) do
17   states :=  $s$ 

```

The replicated service is assumed to be stateful and its implementation realized as a deterministic state machine. That is, starting from the same state und provided with the same command as input, the service implementation, if correct, always yields the same result independent of the instance and independent of the replica that executes it. Listing 1 shows the detailed specification of the service implementation. All possible states of the service's deterministic state machine are represented in the set \mathcal{S} . Each instance of the type *Service* starts in the same initial state $s_0 \in \mathcal{S}$ that is assigned to the instance variable *state_s*. The main interface of a service instance is its public *invoke* method. *invoke* expects a service command *svccmd* of the type \mathcal{O} that encapsulates the invoked service operation together with all required arguments. Moreover, *invoke* has to be provided with the identifier $c \in \mathcal{C}$ of the client that issued the command and a sequence number *invno* assigned by the client to identify the particular invocation. To distinguish different types of sequence numbers, invocation numbers are of the type ${}^i\mathbb{N} \stackrel{\text{def}}{=} \mathbb{N}$. The method *invoke* does not implement the state transition itself. For that purpose, it relies on the internal method *calculateResult* that deterministically maps a specified current state and the invocation information comprising the client, its invocation number, and the service command to a new state and a return value of the type \mathcal{O}' . Being dependent of the concrete service, the method *calculateResult* is not implemented but only declared. *invoke* calls *calculateResult* on the basis of the state referred to by the instance variable *state_s*. When the call returns, it sets *state_s* to the new state and returns with the return value calculated by the command execution. Besides the main method for invoking operations of the realized service, the type *Service* offers methods to create and install snapshots of the service state.

The interface of the invocation handler that is used by clients to issue commands to the replicated service is depicted in Listing 2. The method *startInvocation* takes a service command

Listing 2: Specification of the invocation handler interface.

```
1 module InvocationHandlerc
2   declare input call startInvocation(svccmd ∈  $\mathcal{O}$ )
3   declare output task invocationCompleted(svcret ∈  $\mathcal{O}'$ )
```

of the type \mathcal{O} as input. The respective return value of the command execution is eventually delivered by the output task *invocationCompleted*. Correct clients are only allowed to start a new invocation when a previous invocation was completed, that is, when no other invocation is in progress.

3.2 Time

The system is assumed to be partially synchronous [9]. More precisely, it is assumed that there exist upper bounds for the delivery of messages and how long it takes a process to carry out particular operations but that these upper bounds are not known a priori.

Listing 3: Specification of timers with individual clock speeds.

```
1 module Timers
3 upon init() do
4   timeouts ∈  $\mathbb{N} \cup \{\emptyset\} := \emptyset$ 
6 upon input call schedule(timeout ∈  $\mathbb{N}$ ) do
7   timeouts := timeout
9 upon internal task tick()
10 with
11   timeouts ≠  $\emptyset$  ∧ timeouts > 0
12 do
13   timeouts := timeouts - 1
15 upon output task timerExpired()
16 with
17   timeouts ≠  $\emptyset$  ∧ timeouts = 0
18 do
19   timeouts :=  $\emptyset$ 
21 upon output call isScheduled() → Bool do
22   return timeouts ≠  $\emptyset$ 
24 upon input call cancel() do
25   timeouts :=  $\emptyset$ 
```

In this model, processes have access to local timers that advance at individual speeds. An upper bound for the relative speeds exists but is unknown. Listing 3 presents the specification of a type *Timer* that can signal when a configured amount of time that is relative to each *Timer* instance elapsed. For that purpose, each instance maintains a variable *timeout_s* that stores the remaining number of clock ticks until the timer expires or *None* if no timeout is currently configured. The input method *schedule* can be used to start the timer. It sets *timeout_s* to the specified number of clock ticks. The internal task *tick* regularly decrements the current value of *timeout_s* by one as long as *timeout_s* is greater than zero. Though, it depends on the relative

speed of the timer instance in which interval the task is executed. Once the value of *timeout_s* reaches zero, the output task *timerExpired* signals that the previously configured number of ticks has elapsed. In addition to the basic functionality, the type *Timer* defines methods to determine if a timeout is currently configured and to stop the timer before *timerExpired* is triggered.

3.3 Communication

The network used by processes to communicate with each other by sending messages can be unreliable. It is allowed to reorder, duplicate, delay, and even drop messages. As a consequence, the assumed upper bound for message delivery only holds if a message is sent sufficiently often. Since the upper bound is not known a priori, it is not known in advance how often this is.

The specification of the network and its conditions can be found in Listing 4. Messages (Listing 4a) are modeled as tuples that contain at least one element identifying the message type. A message of the type MSG can be written as $\langle \text{MSG}, [\dots] \rangle$, where $[\dots]$ is a placeholder for a list of optional elements. From the network's perspective, processes act as sinks that receive any kind of message (Listing 4b). The specification of the network itself (Listing 4c) is based on the one given by Castro [5, p. 28]. The state of a network comprises the set of all connected processes stored in nodes_n and the set wire_n that maintains a pair for all messages that are currently in transmission. The first element of such a pair is the message in question and the second is a set of processes supposed to receive the message. To request the network to transmit a message m , the input method *send* is provided. It takes m as argument together with a subset R of the connected processes and notes R in wire_n as intended recipients for m . Delivering messages is the responsibility of the internal task *transmit*. It removes a recipient r of an existing pair $\langle m, R \rangle$ in wire_n and invokes the *receive* method of r passing m as argument. If no recipient for a message m is left, the pair for m in wire_n can be cleaned up by the internal task *discard*. However, as long as the pair for a message m is contained in wire_n , it is possible that the set of recipients is altered by the task *misbehave*. Modeling the unreliable behavior of the network, *misbehave* can add processes to or remove processes from the current set of intended recipients for m . This includes adding new recipients to an empty set or removing all recipients. The only constraint at this point is that the task is fair, that it does not have any bias towards particular messages. The likelihood for being altered has to be the same for all pairs in wire_n to ensure that messages sent infinitely often are also delivered infinitely often and that messages sent finitely often are also delivered finitely often. Note that the network as specified here does not invent completely new messages, which includes that it does not corrupt sent messages. This is actually more restrictive than necessary. All messages in Hybster are authenticated (see below), which ensures their integrity at the protocol level even if the network could deliver corrupted messages.

3.4 Cryptography

Hybster relies on digital signatures and message authentication codes (MACs) to authenticate messages exchanged over the network. Moreover, cryptographic hash functions are used to reduce the amount of transmitted data.

The authenticity of messages is documented by message certificates (Listing 5a). Message certificates are messages themselves, they are tuples comprising a certificate type, a proof of their validity, and further elements if required. Usually, they are directly attached to the mes-

Listing 4: Specification of the communication system shared by all processes.

(a)

```

1 module Message
2   extends Tuple(MessageType, ...)

4   declare output method type()  $\rightarrow$  MessageType

```

(b)

```

1 module MessageSink
2   declare input call receive( $m \in \text{Message}$ )

```

(c)

```

1 module Networkn

3 upon init() do
4    $nodes_n \in \mathcal{P}(\text{MessageSink}) := \text{set of all connected processes}$ 
5    $wire_n \in \mathcal{P}(\text{Message} \times \mathcal{P}(\text{MessageSink})) := \{\}$ 

7 upon input call send( $m \in \text{Message}, R \in \mathcal{P}(\text{MessageSink})$ )
8   asserts
9      $R \subseteq nodes_n$ 
10  do
11    if  $\langle m, R^* \rangle \in wire_n$  then
12       $wire_n := wire_n \setminus \{\langle m, R^* \rangle\} \cup \{\langle m, R^* \cup R \rangle\}$ 
13    else
14       $wire_n := wire_n \cup \{\langle m, R \rangle\}$ 

16 upon internal task transmit( $m \in \text{Message}, r \in \text{MessageSink}$ )
17   with
18      $\exists \langle m, R \rangle \in wire_n : r \in R$ 
19   do
20      $wire_n := wire_n \setminus \{\langle m, R \rangle\} \cup \{\langle m, R \setminus \{r\} \rangle\}$ 
21     output r.receive(m)

23 upon internal task discard( $m \in \text{Message}$ )
24   with
25      $\exists \langle m, \emptyset \rangle \in wire_n$ 
26   do
27      $wire_n := wire_n \setminus \{\langle m, \emptyset \rangle\}$ 

29 upon internal task misbehave( $m \in \text{Message}, R \in \mathcal{P}(\text{MessageSink}), R' \in \mathcal{P}(\text{MessageSink})$ )
30   with
31      $\exists \langle m, R \rangle \in wire_n$ 
32      $R' \subseteq nodes_n$ 
33   do
34      $wire_n := wire_n \setminus \{\langle m, R \rangle\} \cup \{\langle m, R' \rangle\}$ 

```

sage they certify forming a certified message (Listing 5b). The notation for such a message is $\langle \dots \rangle[\textit{certificate}]$, with $[\textit{certificate}]$ being a placeholder for an identifier of the employed certification method. $\langle \dots \rangle_{\sigma_x}$ denotes a message signed by a process x . If a message is certified by a MAC, it is written as $\langle \dots \rangle_{\mu_{x,y,z}}$, where the underlying secret key is shared among the processes x , y , and z . If a process x issues an array of MACs, a so-called authenticator [7] for the processes y and z , the message is annotated with $\langle \dots \rangle_{\alpha_{x:y,z}}$.

Listing 5: Specification of certified messages.

(a)

```
1 module MessageCertificate
2   extends Message

4   declare output method proof() → Proof
```

(b)

```
1 module CertifiedMessage
2   extends Message

4   declare output method certificate() → MessageCertificate
```

Listings 6a and 6b depict the specifications for providers of signatures and MACs, respectively. An instance of the type *SignatureProvider* is initialized with a public/private key pair. It issues digital signatures on the basis of this pair in the method *createSignature* and verifies signatures in *verifySignatures* using the public key connected to the signature in question. Message authentication codes are issued by *MacProvider* instances. MACs are created and verified in the methods *createMac* and *verifyMac* using the shared key that is configured during the initialization of the instance.

The implementation of a cryptographic hash function is provided by the type *DigestProvider* as shown in Listing 6c. This type offers a global method *digest* that returns the hash calculated for a given argument in form of a message similar to the certificates issued by certification providers described above. The hash function is required to be collision resistant such that $\forall m, m': \text{digest}(m) = \text{digest}(m') \leftrightarrow m = m'$ holds with almost absolute certainty during the period in which the message m is relevant and not outdated from perspective of the protocol.

3.5 Faults

Constraints. A process, be it a replica or a client, can fail at any point in time. The behavior of a failed process is undetermined, it can stop the processing of messages completely or it can start to act maliciously by sending messages that are not correct with regard to the protocol specification. Thus, processes can fail arbitrarily, that is, they can be Byzantine. However, this does not apply to all components of replicas. Part of each replica is a trusted subsystem that is assumed to fail only by crashing. That is, independent of the rest of the process, it is assumed that a trusted subsystem either behaves according to the specification or does not carry out any operation, it is assumed that a trusted subsystem never yields a result that is incorrect. Nonetheless, a process is still regarded as either correct or faulty. If some of the components of a process are faulty in any way, the whole process is deemed as faulty.

Furthermore, it is assumed that the cryptographic assumptions hold under all circumstances. Even if an adversary controlled all faulty processes in the system, it is assumed that the collaborating processes would not be able to impersonate correct processes or find collisions for message hashes. Therefore, adversaries are regarded as computationally restricted.

Given a point t in the time \mathcal{T} of some execution of the system, $\mathcal{R}_f(t)$ denotes the set of all replicas that are faulty at t and $\mathcal{C}_f(t)$ the set of all faulty clients. The recovery of processes is not considered. Therefore, it holds: $\forall t, t' \in \mathcal{T}: t \leq t' \rightarrow (\mathcal{R}_f(t) \cup \mathcal{C}_f(t)) \subseteq (\mathcal{R}_f(t') \cup \mathcal{C}_f(t'))$. Further, if a replica or a client is said to be faulty without a reference to time, there is some

Listing 6: Specification of providers for message certificates and cryptographic hashes.

(a)

```

1 module SignatureProviders

3 upon init(pubkey ∈ PublicKey, privkey ∈ PrivateKey) do
4   pubkeys := pubkey
5   privkeys := privkey

7 upon output call createSignature(m) → MessageCertificate do
8   p := signature for m using privkeys
9   return ⟨SIGNATURE, pubkeys, p⟩

11 upon output call verifySignature(⟨SIGNATURE, pubkey, p⟩, m) → Bool do
12   return use pubkey to verify that p is a valid signature for m

```

(b)

```

1 module MacProviders

3 upon init(key ∈ SharedKey) do
4   keys := key

6 upon output call createMac(m) → MessageCertificate do
7   p := calcMac(m)
8   return ⟨MAC, p⟩

10 upon output call verifyMac(⟨MAC, p⟩, m) → Bool do
11   return p = calcMac(m)

13 upon internal call calcMac(m) → Mac do
14   return MAC for m using keys

```

(c)

```

1 module DigestProvider

3 upon global call digest(m) → Message
4   d := calculate digest of m
5   return ⟨DIGEST, d⟩

```

time in the execution in question in which it does not behave according to the specification. As a consequence, the set of faulty replicas is given by $\mathcal{R}_f \stackrel{\text{def}}{=} \bigcup_{t \in \mathcal{T}} \mathcal{R}_f(t)$ and the set of faulty clients by $\mathcal{C}_f \stackrel{\text{def}}{=} \bigcup_{t \in \mathcal{T}} \mathcal{C}_f(t)$. It is assumed that the number of faulty replicas does not exceed a particular f , thus it is assumed that $|\mathcal{R}_f| \leq f < n$ is always true. Opposed to that, there is no required upper bound for the number of faulty clients, that is, $|\mathcal{C}_f| \leq n_{\text{clients}}$. Complementing the sets of faulty processes, the sets for correct replicas and clients are defined as $\mathcal{R}_c \stackrel{\text{def}}{=} \mathcal{R} \setminus \mathcal{R}_f$ and $\mathcal{C}_c \stackrel{\text{def}}{=} \mathcal{C} \setminus \mathcal{C}_f$, respectively.

Quorums. Since processes can fail silently or omit to send messages and since the assumed upper bounds for message delivery and processing times are not known, a correct process cannot know if an expected message from another process p has not arrived yet because p is faulty or because the network or p are slower than currently assumed. Hence, in situations where messages

are not received in time, processes cannot determine if another process p is faulty or not. Even if p provides a correct message in time, it can be Byzantine or fail at a later point in time. As a consequence, state-machine replication protocols do not rely on single replicas to acknowledge protocol information or states but on sufficiently large subsets of the replica group, so-called *quorums* [7]. Formally, a quorum is defined as a set Q such that $Q \subseteq \mathcal{R} \wedge |Q| \geq qs$, with qs being a minimum quorum size that depends on the properties the particular quorum shall provide. The most important quorums are called *intersecting quorums*. If not stated otherwise, these quorums have the fixed minimum size q and fulfill the properties that (1) there is at least one, possibly faulty, replica in the intersection of every two quorums ($2q > n$) and (2) the number of correct replicas in the system suffices to form a quorum ($n \geq q + f$). An important implication of these properties is that each intersecting quorum contains at least one correct replica ($q > f$). Given a number of replica faults f that has to be tolerated, these properties lead to a minimum configuration of $n = 2f + 1$ required replicas with a minimum size for intersecting quorums of $q = f + 1$. Another type of quorums are *acknowledging quorums*, also called weak quorums. These quorums have the property that at least one correct replica is contained. Thus, their minimum size is always $f + 1$, independent of the number of replicas n . Further, in the minimum configuration they have the same size as intersecting quorums.

4 Specification of Hybster

This section describes the properties Hybster provides and presents a formal specification of the protocol. An informal step-by-step description will be part of a future version of this document.

4.1 System Properties

Hybster provides roughly the same properties as PBFT [5], which is a form of linearizability [10] that takes the arbitrary behavior of Byzantine clients into account. From the perspective of each correct client, the replicated system appears to execute all invocations sequentially in an order that respects the real-time relation between the invocations.

The specification of the properties Hybster satisfies is presented in Listing 7 and is based on the one given by Castro [5, p. 30]. It uses the service and invocation handler specifications introduced in Section 3.1.

Initialization. First, a couple of system-wide variables are initialized ($init()_s$). svc_s takes an instance of the service implementation. The set of issued but yet unprocessed invocations is maintained in $pending_s$ and all undelivered responses in $responses_s$. Invocations are triples of a client (\mathcal{C}), an invocation number ($^i\mathbb{N}$), and a service command (\mathcal{O}), thus $\mathcal{I} \stackrel{\text{def}}{=} \mathcal{C} \times ^i\mathbb{N} \times \mathcal{O}$. Invocation numbers are individual to each client and correct clients assign these numbers only once and monotonically increasing. Therefore, invocations can be identified by the pair of a client and an invocation number ($\mathcal{C} \times ^i\mathbb{N}$). A response to an invocation is formed by its identifier and the value returned by its execution (\mathcal{O}'), that is, $\mathcal{Y} \stackrel{\text{def}}{=} \mathcal{C} \times ^i\mathbb{N} \times \mathcal{O}'$. Moreover, the system-wide variable $last_inv_s$ references a vector that stores for each client the invocation number of the last executed invocation. All entries of this vector are initialized with zero and the first valid invocation number for each client is one.

Listing 7: Specification of linearizability with Byzantine clients.

```

1 module ByzantineLinearizabilitys
2   extends InvocationHandler as c

4 upon init()s do
5   svcs := Service()
6   pendings ∈  $\mathcal{P}(\mathcal{I}) := \{\}$ 
7   responsess ∈  $\mathcal{P}(\mathcal{Y}) := \{\}$ 
8   last_invss:  $\mathcal{C} \rightarrow {}^i\mathbb{N}$ , for all c ∈  $\mathcal{C}$  do last_invss(c) := 0

10 upon init()c do
11   invnoc ∈  ${}^i\mathbb{N} := 0$ 
12   svccmdc ∈  $\mathcal{O} \cup \{\emptyset\} := \emptyset$ 
13   isfaultyc ∈ Bool := False

15 upon input call startInvocation(svccmd ∈  $\mathcal{O}$ )c
16   asserts
17     svccmdc =  $\emptyset$ 
18   do
19     invnoc := invnoc + 1
20     svccmdc := svccmd
21     registerInvocation( $\langle c, invno_c, svccmd \rangle$ )

23 upon internal call registerInvocation(inv' =  $\langle c \in \mathcal{C}, \dots \rangle$ )s do
24   pendings := pendings \ {inv ∈ pendings | inv =  $\langle c, \dots \rangle$ } ∪ {inv'}

26 upon internal task processInvocation(inv =  $\langle c \in \mathcal{C}, invno \in {}^i\mathbb{N}, svccmd \in \mathcal{O} \rangle$ )s
27   with
28     inv ∈ pendings ∧ invno > last_invss(c) ∧ preservesFairness(inv)
29   do
30     pendings := pendings \ {inv}
31     svcret := invoke svcs.invoke(c, invno, svccmd)
32     responsess := responsess \ {resp ∈ responsess | resp =  $\langle c, \dots \rangle$ } ∪ { $\langle c, invno, svcret \rangle$ }
33     last_invss(c) := invno

35 declare internal call preservesFairness(inv =  $\langle c \in \mathcal{C}, invno \in {}^i\mathbb{N}, svccmd \in \mathcal{O} \rangle$ )s

37 upon output task invocationCompleted(svcret ∈  $\mathcal{O}'$ )c
38   with
39     isfaultyc ∨ resp =  $\langle c, invno_c, svcret \rangle$  ∈ responsess
40   do
41     responsess := responsess \ {resp}
42     svccmdc :=  $\emptyset$ 

44 upon internal task clientFailureOccured()c
45   with
46     ¬isfaultyc
47   do
48     isfaultyc := True

50 upon internal task injectFaultyInvocation(c ∈  $\mathcal{C}, invno \in {}^i\mathbb{N}, svccmd \in \mathcal{O}$ )c
51   with
52     isfaultyc
53   do
54     registerInvocation( $\langle c, invno, svccmd \rangle$ )

```

Invocation. Invocation handler instances are dedicated to single clients and used by them to invoke operations of the replicated service. A handler instance maintains three variables (see *init()*_{*c*}), the invocation number of the last invocation *invno*_{*c*}, *svccmd*_{*c*} that stores the service command of a currently running but not yet completed invocation or \emptyset if no invocation is in progress, and the flag *isfaulty*_{*c*} that marks faulty clients. Correct clients start only one invocation at a time and wait for its completion before they start a new one. Hence, the input method *startInvocation* expects that *svccmd*_{*c*} is \emptyset when it is called. Subsequently, it assigns the next invocation number to the passed service command, sets *svccmd*_{*c*} to mark the invocation as in progress, and submits the new invocation by calling *registerInvocation*. This system-wide method adds invocations to the set of pending invocations while ensuring that each client has at most one outstanding invocation in that set.

Processing. Pending invocations are delivered to the service instance by the task *processInvocation*, one invocation at a time. It considers only invocations with numbers higher than the numbers stored in *last_inv*_{*s*}, that is, with numbers higher than the one of the last invocation that was executed for the respective client. Being always the case for correct clients, this prevents faulty clients from injecting invocations in a non-increasing order. Further, *processInvocation* makes use of the method *preservesFairness* to select the invocation to be executed next. This auxiliary method is only declared and documents the requirement that the selection mechanism is unbiased and does not give particular clients an advantage. Once an invocation has been chosen, it is removed from the pending set and handed over to the service implementation. The value returned by the service is used to create a response. Before this response is added to the set of undelivered responses, older responses for the client are removed if necessary. A faulty client could have sent a new invocation without collecting the response for a previous one. Finally, *last_inv*_{*s*} is updated to mark the current invocation as executed.

Completion. Adhering to the interface specified by the type *InvocationHandler*, invocations are delivered to clients through the output task *invocationCompleted*. Provided that a client *c* is not marked as faulty, this task removes a response belonging to *c* from the set of undelivered responses and marks the invocation as completed by setting *svccmd*_{*c*} to *None* before it delivers the return value.

Faulty Clients. A client *c* can fail at any point in time. This is modeled by the task *clientFailureOccured*. If executed, it simply sets the flag *isfaulty*_{*c*} to *True*. From that time on, the client *c* can circumvent the regular *startInvocation* method and can issue arbitrary invocation as realized by the task *injectFaultyInvocation*. However, it can still only do so for invocations that are attributed to *c* itself. As stated in Section 3.5, the system model assumes that even arbitrarily behaving clients are not able to impersonate correct clients. In addition to the injection of arbitrary invocation, a faulty client can locally deliver any return value as reflected in the task *invocationCompleted*.

4.2 Auxiliary Modules

The specifications of Hybster's main protocols make use of two mechanisms that help to cope with the partially synchronous environment and the unreliable network as assumed by the system model (Section 3): adaptive timeouts and retransmitting connections. The respective modules realizing these mechanisms are presented in the following.

4.2.1 Adaptive Timeouts

The time model Hybster's specification is based on assumes upper bounds for message delivery and processing times, upper bounds, however, that are unknown (cf. Section 3.2). As a consequence, it cannot be known how long certain protocol actions need to finish. Nonetheless, it is known that they do not require an infinite amount of time as long as their completion does not depend on a faulty process. Therefore, provided that some protocol action can be retried infinitely often, it is possible to approach the upper bound by increasing the value for the timeout that determines when the action is regarded as failed.

Listing 8: Specification of timeouts adaptive to unknown upper bounds for delays.

```

1 module AdaptiveTimeouts

3 upon init(procs ∈ Set) do
4   procss := procs
5   timers := Timer()
6   delay_deltas ∈ ℕ := increment if timeout expired
7   delays ∈ ℕ := delay_deltas
8   expireds ∈ Bool := False

10 upon input call start(x)
11   asserts
12     x ∈ procss // x could be used in more sophisticated implementations
13   do
14     adapt()
15     output timers.schedule(delays)

17 upon input call adapt() do
18   if expireds then
19     delays := delays + delay_deltas
20     expireds := False

22 upon input call stop() do
23   output timers.cancel()
24   expireds := False

26 upon internal task timerExpired()
27   with
28     input timers.timerExpired()
29   do
30     expireds := True

32 upon output call isExpired() → Bool do
33   return expireds

35 upon output call isRunning() → Bool do
36   return input timers.isScheduled()

```

This mechanism is realized in the module *AdaptiveTimeout* shown in Listing 8. It is intended that one instance of this module is used for one particular protocol action that has to be carried out by one process from a group of processes. For example, if a replica is supposed to provide a particular message, an instance of this module can be used to signal when this step in the protocol is considered as failed. If this is the case, the step is repeated, but this time, a different replica is chosen to provide the message and it is given more time to do so. Eventually, this

step will complete successfully, at the latest when a correct replica is selected and the currently configured timeout exceeded the upper bound for providing the expected message.

Although the implementation of *AdaptiveTimeout* is tailored to the time model stated above, please note that *AdaptiveTimeout* in fact abstracts from the concrete model. Currently, it increases the timeout value linearly. If a time model as in PBFT were used in which a delay function $delay(t)$ does not grow faster than the time t [5], an exponential increase could compensate for that. More sophisticated strategies would be required if the time bounds were only assumed to hold long enough such that the system is able to make progress [9]. Still, higher-level protocols could rely on the abstraction provided by *AdaptiveTimeout*.

The module *AdaptiveTimeout* in detail: Each instance maintains a reference to a timer instance, the current timeout value stored in $delay_s$ and initialized with a configured constant $delay_delta_s$, and a flag $expired_s$ set when the timeout expired. The instance also gets the set of processes that are within the considered group. This set is actually not required by the presented implementation and is only provided to allow for other strategies that adjust the timeout value for processes individually instead of globally for all processes. The timeout is started through the input method *start*. It takes the process as argument that is currently expected to carry out the monitored action. Again, the presented implementation does not make use of this information. Upon being invoked, the method *start* configures the timer to signal after time $delay_s$. Prior to that, it is checked if the current timeout value $delay_s$ needs to be adjusted. The strategy for that is realized by the input method *adapt*. If the timer was marked as expired, it adds $delay_delta_s$ to $delay_s$ and resets the expired flag. If the timed protocol action completed or if the timeout is not needed anymore, the input method *stop* can be invoked canceling a currently running timer. Otherwise, the timer will eventually expire, monitored with the internal *timerExpired* task. This sets the expired flag which can be queried by the output method *isExpired*. Additionally, *isRunning* returns if the timer is currently scheduled or not.

4.2.2 Connections

Even if a process expected to provide some message were correct and even if other processes waited actually long enough to receive this message, the message might not come in time since the network could have dropped it. The upper bounds for message delivery only hold if the message is sent sufficiently, virtually infinitely often (see Section 3.3).

Retransmitting messages is the responsibility of the *Connection* module presented in Listing 9. It sends registered messages over and over again until they are explicitly removed from the connection by a higher-level protocol. An instance of the *Connection* module expects a set of recipients when constructed. It further has a reference to the network and keeps with out_s a set of all messages that are currently in transmission. Messages are handed over to the connection via the *send* input method which adds the message to out_s . The internal task *retransmit* is enabled as long as some message is contained in out_s . When executed, it picks one message and (re-)sends it to the configured recipients. Nevertheless, the selected message remains in the set out_s and could be chosen again in a later execution of the task. If multiple messages are contained in out_s , the selection mechanism is supposed to be fair, on average all messages are retransmitted at the same rate. The module offers multiple variants of the method *remove* used to deregister messages from the connection. *remove* can be invoked with a single message, a set of messages, or a predicate that determines which messages are to be removed from out_s .

Listing 9: Specification of a retransmitting connection masking message loss.

```

1 module Connections

3 upon init(recipients ∈  $\mathcal{P}(\text{MessageSink})$ ) do
4   nets ∈ Network := reference to the network
5   recipientss := recipients
6   outs ∈  $\mathcal{P}(\text{Message})$  := {}

8 upon input call send(m ∈ Message) do
9   outs := outs ∪ {m}

11 upon internal task retransmit(m ∈ Message)
12   with
13     m ∈ outs
14   do
15     output nets.send(m, recipientss)

17 upon input call remove(m ∈ Message) do
18   outs := outs \ {m}

20 upon input call remove(M ∈  $\mathcal{P}(\text{Message})$ ) do
21   outs := outs \ M

23 upon input call remove(P: Message → Bool) do
24   outs := outs \ {m ∈ outs | P(m)}

```

4.3 TrInX

Listing 10: Specification of the trusted subsystem TrInX.

```

1 module TrInXs

3 upon init(id ∈ TssID, ctrtypes ∈ Set, key ∈ SharedKey) do
4   insts := id
5   ctrtypess := ctrtypes
6   macprovs := MacProvider(key)
7   ctrs: ctrtypes → ℕ, for all tc ∈ ctrtypes do ctrs(tc) := 0

9 upon public call createContinuingCertificate(tc, tv' ∈ ℕ, m) → MessageCertificate
10  asserts
11    tc ∈ ctrtypess ∧ tv' ≥ ctrs(tc)
12  do
13    tv := ctrs(tc)
14    p := macprovs.createMac(insts||tc||tv'||tv||m)
15    ctrs(tc) := tv'
16    return ⟨TCTR, insts, tc, tv', tv, p⟩

18 upon public call createIndependentCertificate(tc, tv' ∈ ℕ, m) → MessageCertificate
19  asserts
20    tc ∈ ctrtypess ∧ tv' > ctrs(tc)
21  do
22    p := macprovs.createMac(insts||tc||tv'|| − ||m)
23    ctrs(tc) := tv'
24    return ⟨TCTR, insts, tc, tv', −, p⟩

26 upon output call verifyCertificate(⟨TCTR, inst, tc, tv', tv, p⟩, m) → Bool do
27   return macprovs.verifyMac(inst||tc||tv'||tv||m)

29 upon input call forwardCounter(tc, tv' ∈ ℕ)
30  asserts
31    tc ∈ ctrtypess ∧ tv' ≥ ctrs(tc)
32  do
33    ctrs(tc) := tv'

```

4.4 Processes

4.4.1 Clients

Listing 11: Specification of invocation handlers dedicated to a client.

```
1 module ByzInvocationHandlerc
2   extends InvocationHandler
3   extends MessageSink

5 upon init(pubkey ∈ PublicKey, privkey ∈ PrivateKey, pubkeys:  $\mathcal{R} \rightarrow \text{PublicKey}$ ) do
6   repconnsc:  $\mathcal{R} \rightarrow \text{Connection}$ , for all  $r \in \mathcal{R}$  do repconnsc( $r$ ) := Connection( $\{r\}$ )
7   repgrpconnc := Connection( $\mathcal{R}$ )
8   sigprovc := SignatureProvider(pubkey, privkey)
9   pubkeysc := pubkeys
10  imc := InvocationHandlerc::Invocation()

12 upon input call startInvocation(svcCmd ∈  $\mathcal{O}$ ) do
13   output imc.startInvocation(svcCmd)

15 upon output task invocationCompleted(svcret ∈  $\mathcal{O}'$ ) with
16   input imc.invocationCompleted(svcret)

18 upon input call receive( $m \in \text{Message}$ ) do
19   output imc.receive( $m$ )

21 upon internal call hasValidCertificate( $m = \langle \dots \rangle_{\sigma_r} \rightarrow \text{Bool}$ ) do
22    $p := m.\text{certificate}().\text{proof}()$ 
23   return input sigprovc.verifySignature( $\langle \text{SIGNATURE}, \text{pubkeys}_c(r), p \rangle, m$ )
```

4.4.2 Replicas

Listing 12: Specification of the current view status used by protocol modules.

```
1 module ViewStatuss

3 upon init(replica ∈  $\mathcal{R}$ ) do
4   replicas ∈  $\mathcal{R} := \text{replica}$ 
5   v_stabs ∈  ${}^v\mathbb{N} := 0$ 
6   v_curs ∈  ${}^v\mathbb{N} := v\_stab_s$ 

8 upon global call leader( $v \in {}^v\mathbb{N}$ ) →  $\mathcal{R}$  do
9   return  $r_{v \bmod n}$ 

11 upon input call leaveViewFor( $v' \in {}^v\mathbb{N}$ )
12   asserts
13      $v\_cur_r < v'$ 
14   do
15      $v\_cur_r := v'$ 

17 upon input call enterView( $v' \in {}^v\mathbb{N}$ )
18   asserts
19      $v\_stab_r < v' \leq v\_cur_r$ 
20   do
21      $v\_stab_r := v'$ 

23 upon output call latestLeader() →  $\mathcal{R}$  do
24   return leader( $v\_stab_s$ )
```



```

26 upon output call isInStableView()  $\rightarrow \text{Bool}$  do
27   return  $v\_cur_s = v\_stab_s$ 

29 upon output call isStableLeader()  $\rightarrow \text{Bool}$  do
30   return  $isInStableView() \wedge replica_s = latestLeader()$ 

32 upon output call isStableFollower()  $\rightarrow \text{Bool}$  do
33   return  $isInStableView() \wedge replica_s \neq latestLeader()$ 

35 upon output call stableView()  $\rightarrow {}^v\mathbb{N}$  do
36   return  $v\_stab_s$ 

38 upon output call currentView()  $\rightarrow {}^v\mathbb{N}$  do
39   return  $v\_cur_s$ 

```

Listing 13: Specification of a replica making use of multiple protocol modules.

```

1 module Replicar
2   extends MessageSink

4 def  $\mathcal{TC} \stackrel{\text{def}}{=} \{\mathfrak{M}, \mathfrak{D}, \mathfrak{N}\}$ 

6 upon init( $pubkey \in \text{PublicKey}, privkey \in \text{PrivateKey},$ 
    $tsskey \in \text{SharedKey}, pubkeys: \mathcal{R} \cup \mathcal{C} \rightarrow \text{PublicKey}, tssinsts: \mathcal{R} \rightarrow \text{TssID}$ ) do
7    $repconns_r: \mathcal{R} \setminus \{r\} \rightarrow \text{Connection}, \text{for all } r \in \mathcal{R} \setminus \{r\} \text{ do } repconns_r(r) := \text{Connection}(\{r\})$ 
8    $repgrpconn_r := \text{Connection}(\mathcal{R} \setminus \{r\})$ 
9    $tss_r := \text{TrInX}(tssinsts(r), \mathcal{TC}, tsskey)$ 
10   $tssinsts_r := tssinsts$ 
11   $sigprov_r := \text{SignatureProvider}(pubkey, privkey)$ 
12   $pubkeys_r := pubkeys$ 
13   $cm_r := \text{Replica}_r::\text{Client}$ 
14   $om_r := \text{Replica}_r::\text{Ordering}$ 
15   $km_r := \text{Replica}_r::\text{Checkpointing}$ 
16   $vm_r := \text{Replica}_r::\text{ViewChange}$ 
17   $em_r := \text{Replica}_r::\text{Execution}$ 
18   $sm_r := \text{Replica}_r::\text{StateTransfer}$ 

20 upon input call receive( $m \in \text{Message}$ ) do
21   // dispatch message according to message type
22   output  $(cm_r | om_r | km_r | vm_r | em_r | sm_r).receive(m)$ 

24 upon internal call hasValidCertificate( $m = \langle \dots \rangle_{\sigma_x} \rightarrow \text{Bool}$ ) do
25    $p := m.certificate().proof()$ 
26   return input  $sigprov_c.verifySignature(\langle \text{SIGNATURE}, pubkeys_r(x), p \rangle, m)$ 

28 upon internal call hasValidCertificate( $m = \langle \dots \rangle_{\tau(r', tc, tv', tv)} \rightarrow \text{Bool}$ ) do
29    $p := m.certificate().proof()$ 
30   return input  $tss_r.verifyCertificate(\langle \text{TCTR}, tssinsts_r(r'), tc, tv', tv, p \rangle, m)$ 

```

4.5 Invocation Protocol

4.5.1 Client Side

Listing 14: Specification of the client role of the invocation protocol.

```

1 module InvocationHandlerc::Invocation

3 upon internal call leader( $v \in {}^v\mathbb{N}$ )  $\rightarrow \mathcal{R}$  do
4   return  $r_{v \bmod n}$ 

6 upon init() do
7    $v_c \in {}^v\mathbb{N} := 0$ 
8    $invno_c \in {}^i\mathbb{N} := 0$ 
9    $req_c \in Message \cup \{\emptyset\} := \emptyset$ 
10   $reps_c \in \mathcal{P}(Message) := \{\}$ 
11   $rq\_timeout_c := AdaptiveTimeout(\mathcal{R})$ 
12   $retrans_c \in Bool := False$ 

14 upon input call startInvocation( $svccmd \in \mathcal{O}$ )
15  asserts
16     $req_c = \emptyset$ 
17  do
18     $invno_c := invno_c + 1$ 
19     $req_c := \langle REQUEST, c, invno_c, svccmd \rangle_{\sigma_c}$ 
20     $sigprov_c.createSignature(req_c)$ 
21    output  $repconn_s_c(leader(v_c)).send(req_c)$ 
22     $rq\_timeout_c.start(leader(v_c))$ 

24 upon internal task retransmit()
25  with
26     $req_c \neq \emptyset \wedge \neg retrans_c \wedge rq\_timeout_c.isExpired()$ 
27  do
28     $rq\_timeout_c.adapt()$ 
29    output  $repgrpconn_c.send(req_c)$ 
30     $retrans_c := True$ 

32 upon global call isCorrectReply( $rp, r \in \mathcal{R}, v \in {}^v\mathbb{N}, c' \in \mathcal{C}, invno \in {}^i\mathbb{N}, svcret \in \mathcal{O}'$ ) do
33  return  $rp = \langle REPLY, r, v, c', invno, svcret \rangle_{\sigma_r} \wedge hasValidCertificate(rp)$ 

35 upon input call receive( $rp' = \langle REPLY, r, v', \cdot, \cdot, svcret' \rangle_{\dots}$ ) do
36  if  $req_c \neq \emptyset \wedge isCorrectReply(rp', \cdot, \cdot, c, invno_c, \cdot)$  then
37    if  $\nexists v \nexists svcret: rp = \langle REPLY, r, v, c, invno_c, svcret \rangle_{\dots} \in reps_c$  then
38       $reps_c := reps_c \cup \{rp'\}$ 
39    else if  $svcret' \neq svcret$  then
40      // r is faulty
41    else if  $v' > v$  then
42       $reps_c := reps_c \cup \{rp'\} \setminus \{rp\}$ 

44 upon global call isCorrectReplyCertificate( $Y \subseteq Message, c' \in \mathcal{C}, invno \in {}^i\mathbb{N}, svcret \in \mathcal{O}'$ ) do
45  return  $|Y| > f$ 
46     $\wedge (\forall rp \in Y \exists r \in \mathcal{R}: isCorrectReply(rp, r, \cdot, c', invno, svcret))$ 
47     $\wedge (\forall rp' \in Y: rp' = \langle REPLY, r, \dots \rangle_{\dots} \rightarrow rp' = rp)$ 

49 upon output task invocationCompleted( $svcret \in \mathcal{O}'$ )
50  with
51     $\exists Y \subseteq reps_c: isCorrectReplyCertificate(Y, c, invno_c, svcret)$ 
52  do

```

```

53  rq_timeoutc.stop()
54  output repconnsc(leader(vc)).remove(reqc)
55  if retransc then
56    output repgrpconnc.remove(reqc)
57    vc := max({v | (REPLY, ·, v, ...) ∈ Y})
58    reqc := ∅
59    repsc := {}
60    retransc := False

```

4.5.2 Replica Side

Listing 15: Specification of the replica role of the invocation protocol.

```

1  module Replicar::Client

3  def global ReturnEntry  $\stackrel{\text{def}}{=} {}^i\mathbb{N} \times (\mathcal{O}' \cup \{\emptyset\})$ 

5  upon global call invocationNumber((invno, svcret) ∈ ReturnEntry) →  ${}^i\mathbb{N}$  do
6    return invno

8  upon init() do
9    cstatusr := ViewStatus(r)
10   reqsr ∈  $\mathcal{P}(\text{Message})$  := {}
11   rq_timeoutsr:  $\mathcal{C} \rightarrow \text{AdaptiveTimeout}$ 
12   last_retvalsr:  $\mathcal{C} \rightarrow \text{ReturnEntry}$ 
13   cliconnsr:  $\mathcal{C} \rightarrow \text{Connection}$ 
14   for all c ∈  $\mathcal{C}$  do
15     rq_timeoutsr(c) := AdaptiveTimeout( $\mathcal{R}$ )
16     last_retvalsr(c) := (0, ∅)
17     cliconnsr(c) := Connection({c})
18   contactr := cstatusr.latestLeader()

20 upon global call isCorrectRequest(rq, c ∈  $\mathcal{C}$ , invno ∈  ${}^i\mathbb{N}$ , svccmd ∈  $\mathcal{O}$ ) do
21   return rq = (REQUEST, c, invno, svccmd)σc ∧ hasValidCertificate(rq)

23 upon input call receive(rq' = (REQUEST, c, invno', svccmd)...) do
24   if isRelevantRequest(rq') ∧ isCorrectRequest(rq', ...) then
25     storeRequest(rq')
26     if cstatusr.isStableFollower() then
27       forwardRequest(rq')
28     startRequestTimeout(c)

30 upon internal call isRelevantRequest(rq' = (REQUEST, c, invno', ·)...) do
31   return ¬isRequestExecuted(rq') ∧ (∄ rq ∈ reqsr ∃ invno ≥ invno': rq = (REQUEST, c, invno, ·)...)

33 upon internal call storeRequest(rq' = (REQUEST, c, ...)...) do
34   reqsr := reqsr \ {rq ∈ reqsr | rq = (REQUEST, c, ...)...)
35   reqsr := reqsr ∪ {rq'}

37 upon internal call forwardRequest(rq = (REQUEST, c, ·)...) do
38   output repconnsr(contactr).remove({m | m = (REQUEST, c, ...)...})
39   output repconnsr(contactr).send(rq)

41 upon output call pendingRequests() →  $\mathcal{P}(\text{Message})$  do
42   return reqsr

44 upon input call requestExecuted(rq' = (REQUEST, c, invno', svccmd)..., svcret ∈  $\mathcal{O}'$ ) do

```

```

45  reqsr := reqsr \ {rq | rq = ⟨REQUEST, c, invno, ·⟩... ∧ invno ≤ invno'}
46  last_retvalsr(c) := ⟨invno', svcret⟩
47  sendReply(c, invno', svcret)
48  if  $\nexists$ ⟨REQUEST, c, ...⟩... ∈ reqsr ∧ rq_timeoutsr(c).isRunning() then
49    rq_timeoutsr(c).stop()

51 upon output call isRequestExecuted(rq = ⟨REQUEST, c, invno, ·⟩...) do
52   return invno ≤ invocationNumber(last_retvalsr(c))

54 upon internal call sendReply(c ∈ C, invno ∈ iN, svcret ∈ O') do
55   rp := ⟨REPLY, r, cstatusr.stableView(), c, invno, svcret⟩σr
56   sigprovr.createSignature(rp)
57   output cliconnsr(c).remove({m | m = ⟨REPLY, r, ·, c, ...⟩...})
58   output cliconnsr(c).send(rp)

60 upon output call createReturnValueMapSnapshot() → (C → ReturnEntry) do
61   return last_retvalsr

63 upon input call installReturnValueMap(retvals ∈ C → ReturnEntry) do
64   last_retvalsr := retvals
65   for all rq = ⟨REQUEST, c, ...⟩... ∈ reqsr do
66     if isRequestExecuted(rq) then
67       reqsr := reqsr \ {rq}
68     if rq_timeoutsr(c).isRunning() then
69       rq_timeoutsr(c).stop()

71 upon internal call startRequestTimeout(c ∈ C) do
72   if rq_timeoutsr(c).isRunning() then
73     rq_timeoutsr(c).stop()
74   rq_timeoutsr(c).start(contactr)

76 upon output call suspectsLeader() do
77   return  $\exists to \in rq\_timeouts_r : to.isExpired()$ 

79 upon input call leaveViewFor(v' ∈ vN) do
80   cstatus.leaveViewFor(v')

82 upon input call enterView(v' ∈ vN) do
83   cstatus.enterView(v')
84   if cstatus.isInStableView then
85     output repconnsr(contactr).remove({rq | rq = ⟨REQUEST, ...⟩...})
86     contactr := cstatus.latestLeader()
87     for all rq = ⟨REQUEST, c, ...⟩... ∈ reqsr do
88       if cstatus.isStableFollower() then
89         forwardRequest(rq)
90         startRequestTimeout(c)

```

4.6 Ordering Protocol

Listing 16: Specification of the ordering protocol.

```

1 module Replicar::Ordering

3 uses vmr ∈ Replicar::ViewChange
4 uses cmr ∈ Replicar::Client
5 uses emr ∈ Replicar::Execution

7 // An abort state comprises a prepared set P and
8 // the value of the trusted counter  $\mathfrak{D}$  at the time of the abort.
9 def global AbortState  $\stackrel{\text{def}}{=} (P \in \mathcal{P}(\text{Message})) \times \mathbb{N}$ 

11 let order_wnd ∈  $\mathbb{N}_1$ 

13 upon global call orderWindowSize() do
14   return order_wnd

16 upon init() do
17   ostatusr := ViewStatus(r)
18   omgsr ∈  $\mathcal{P}(\text{Message}) = \{\}$ 
19   o_actr ∈  ${}^o\mathbb{N} := 0$ 
20   o_baser ∈  ${}^o\mathbb{N} := 0$ 
21   o_commr ∈  ${}^o\mathbb{N} := o\_base_r$ 
22   o_prepr ∈  ${}^o\mathbb{N} := o\_comm_r$ 
23   o_maxr  $\stackrel{\text{def}}{=} o\_base_r + order\_wnd$ 

25 upon internal task proposalReady(rq = ⟨REQUEST, c, invno, svccmd⟩σc)
26   with
27     o' = o_prepr + 1
28     ostatusr.isStableLeader() ∧ o' ≤ o_maxr
29     v = ostatusr.stableView()
30     rq ∈ input cmr.pendingRequests() ∧ (∄o: ⟨PREPARE, r, v, o, rq⟩... ∈ omgsr)
31   do
32     sendPrepare(v, o', rq)
33     o_prepr := o'
34     o_actr := o'

36 upon internal call sendPrepare(v ∈  ${}^v\mathbb{N}$ , o ∈  ${}^o\mathbb{N}$ , rq ∈ Message) do
37   pr := ⟨PREPARE, r, v, o, rq⟩τ(r,  $\mathfrak{D}$ , v|o, -)
38   invoke tssr.createIndependentCertificate( $\mathfrak{D}$ , v|o, pr)
39   omgsr := omgsr ∪ {pr}
40   output repgrpconnr.send(pr)

42 upon global call isCorrectPrepare(pr, l ∈  $\mathcal{R}$ , v ∈  ${}^v\mathbb{N}$ , o ∈  ${}^o\mathbb{N}$ , rq ∈ Message) do
43   return l = leader(v) ∧ pr = ⟨PREPARE, l, v, o, rq⟩τ(l,  $\mathfrak{D}$ , v|o, -) ∧ hasValidCertificate(pr)
     ∧ isCorrectRequest(rq, ...)

45 upon input call receive(pr = ⟨PREPARE, l, v, o, rq⟩...) do
46   if l ≠ r ∧ v = ostatusr.stableView() ∧ inOrderWindow(o)
     ∧ isCorrectPrepare(pr, ...) then
47     omgsr := omgsr ∪ {pr}

49 upon internal task nextPrepared(pr = ⟨PREPARE, l, v, o', rq⟩...)
50   with
51     ostatusr.isStableFollower()

```

```

52    $pr \in msgs_r \wedge o' = o\_prep_r + 1$ 
53   do
54      $o\_prep_r := o'$ 
55      $sendCommit(v, o', digest(rq))$ 
56      $o\_act_r := o'$ 

58 upon internal call  $sendCommit(v \in {}^v\mathbb{N}, o \in {}^o\mathbb{N}, d \in \langle DIGEST \dots \rangle)$  do
59    $co := \langle COMMIT, r, v, o, d \rangle_{\tau(r, \mathfrak{D}, v|o, -)}$ 
60   invoke  $tss_r.createIndependentCertificate(\mathfrak{D}, v|o, co)$ 
61    $msgs_r := msgs_r \cup \{co\}$ 
62   output  $repgrpconn_r.send(co)$ 

64 upon global call  $isCorrectCommit(co, r' \in \mathcal{R}, v \in {}^v\mathbb{N}, o \in {}^o\mathbb{N}, d \in \langle DIGEST \dots \rangle)$  do
65   return  $r' \neq leader(v) \wedge co = \langle COMMIT, r', v, o, d \rangle_{\tau(r', \mathfrak{D}, v|o, -)} \wedge isValidCertificate(co)$ 

67 upon input call  $receive(co = \langle COMMIT, r', v, o, d \rangle \dots)$  do
68   if  $r' \neq r \wedge v = ostatus_r.stableView() \wedge \neg isCommitted(o) \wedge inOrderWindow(o)$   

    $\wedge isCorrectCommit(co, \dots)$  then
69      $msgs_r := msgs_r \cup \{co\}$ 

71 upon global call  $isCorrectCommittedCertificate(O \subseteq Message, v \in {}^v\mathbb{N}, o \in {}^o\mathbb{N}, rq \in Message)$  do
72   return  $|O| \geq q \wedge (\exists pr \in O: isCorrectPrepare(pr, \cdot, v, o, rq)$   

    $\wedge (\forall m \in O: m = pr \vee isCorrectCommit(m, \cdot, v, o, digest(rq)))$ 

74 upon internal task  $nextCommitted(rq = \langle REQUEST, c, t, svccmd \rangle \dots)$ 
75   with
76      $o' = o\_comm_r + 1$ 
77      $\exists O \subseteq msgs_r: isCorrectCommittedCertificate(O, \cdot, o', rq)$ 
78   do
79      $markCommitted(o')$ 
80     output  $em_r.executeRequest(o', rq)$ 

82 upon input call  $forwardOrderWindow(o' \in {}^o\mathbb{N})$ 
83   asserts
84      $o' > o\_base_r$ 
85   do
86      $o\_base_r := o'$ 
87      $discardOrderMessages()$ 
88     if  $\neg isCommitted(o')$  then
89        $markCommitted(o')$ 

91 upon internal call  $discardOrderMessages()$  do
92    $msgs_r := \{m \in msgs_r \mid inOrderWindow(orderNumber(m))\}$ 
93   output  $repgrpconn_r.remove(\{m \mid (m = \langle PREPARE, \dots \rangle \dots \vee m = \langle COMMIT, \dots \rangle \dots) \wedge m \notin msgs_r\})$ 

95 upon internal call  $markCommitted(o \in {}^o\mathbb{N})$ 
96   asserts
97      $o > o\_comm_r$ 
98   do
99     if  $o > o\_prep_r$  then
100        $o\_prep_r := o$ 
101        $o\_comm_r := o$ 

103 upon output call  $isCommitted(o \in {}^o\mathbb{N})$  do
104   return  $o \leq o\_comm_r$ 

106 upon output call  $lastCommitted()$  do

```

```

107 return  $o\_comm_r$ 

109 upon output call  $inOrderWindow(o \in {}^o\mathbb{N})$  do
110   return  $o\_base_r < o \leq o\_max_r$ 

112 upon output call  $orderWindowBase()$ 
113   return  $o\_base_r$ 

115 upon output call  $obtainOrderingState(o\_base \in {}^o\mathbb{N})$  do
116   return  $\{m \in msgs_r \mid orderNumber(m) > o\_base\}$ 

118 upon global call  $isCorrectOrderingState(I \subseteq Message, v \in {}^v\mathbb{N}, o\_base \in {}^o\mathbb{N})$  do
119   return  $\forall m \in I \exists o: (isCorrectPrepare(m, \cdot, v, o, \cdot) \vee isCorrectCommit(m, \cdot, v, o, \cdot))$ 
       $\wedge o\_base < o \leq o\_base + order\_wnd$ 

121 upon input call  $installOrderingState(I \subseteq Message)$ 
122   asserts
123    $isCorrectOrderingState(I, \cdot, ostatus_r.stableView())$ 
124   do
125      $msgs_r := msgs_r \cup \{m \in I \mid inOrderWindow(orderNumber(m))\}$ 

127 upon output call  $createAbortState() \rightarrow AbortState$  do
128    $P := \{\langle PREPARE, \cdot, \cdot, o, \cdot \rangle \dots \in msgs_r \mid o - 1 = o\_base_r \vee \langle PREPARE, \cdot, \cdot, o - 1, \cdot \rangle \dots \in msgs_r\}$ 
129   return  $\langle P, ostatus_r.currentView() \mid o\_act_r \rangle$ 

131 upon input call  $abortView(v' \in {}^v\mathbb{N}, m \in \langle \dots \rangle_{\tau(\dots)})$  do
132   invoke  $tss_r.createConsecutiveCertificate(\mathfrak{D}, v' | 0, m)$ 
133    $o\_act_r := 0$ 
134    $ostatus_r.leaveViewFor(v')$ 

136 upon global call  $isCorrectAbortState($ 
       $b = \langle P, tv\_last \rangle \in AbortState, v \in {}^v\mathbb{N}, v' \in {}^v\mathbb{N}, o\_base \in {}^o\mathbb{N}, m \in \langle \dots \rangle_{\tau(r', \mathfrak{D}, v' | 0, tv\_last)})$  do
137   return  $hasValidCertificate(m) \wedge isCorrectAbortState(b, v, v', o\_base)$ 

139 upon global call  $isCorrectAbortState(\langle P, tv\_last \rangle \in AbortState, v \in {}^v\mathbb{N}, v' \in {}^v\mathbb{N}, o\_base \in {}^o\mathbb{N})$  do
140   return  $tv\_last = \langle v, o\_ctr \rangle$ 
141    $\wedge (\exists o\_max: isCorrectPrepareSet(P, v, o\_base, o\_max) \wedge o\_max \geq o\_ctr)$ 

143 upon global call  $isCorrectPrepareSet(P, v \in {}^v\mathbb{N}, o\_base \in {}^o\mathbb{N}, o\_max \in {}^o\mathbb{N})$  do
144   return  $|P| = 0 \wedge o\_max = o\_base \vee 0 < |P| \leq order\_wnd$ 
145    $\wedge o\_base = minOrderNumber(P) - 1 \wedge o\_max = maxOrderNumber(P)$ 
146    $\wedge (\forall pr \in P \exists o: isCorrectPrepare(pr, v, o, \cdot))$ 
147    $\wedge (\forall o' \in {}^o\mathbb{N}: o\_base < o' \leq o\_max \rightarrow \langle PREPARE, \cdot, \cdot, o' \rangle \dots \in P)$ 

149 upon input call  $abortView(v' \in {}^v\mathbb{N})$  do
150   invoke  $tss_r.forwardCounter(\mathfrak{D}, v' | 0)$ 
151    $o\_act_r := 0$ 
152    $ostatus_r.leaveViewFor(v')$ 

154 upon input call  $installAbortState(b = \langle P, tv\_last \rangle \in AbortState)$ 
155   asserts
156    $isCorrectAbortState(b, ostatus_r.stableView(), \dots)$ 
157   do
158      $installOrderingState(P)$ 

160 upon output call  $createViewTransition(v' \in {}^v\mathbb{N}, B \subseteq AbortState)$  do
161    $P^* := \{pr \mid \langle P, \cdot \rangle \in B \wedge pr \in P\}$ 

```

```

162 if  $|P^*| = 0$  then
163   return  $\emptyset$ 
164 else
165    $o\_max := \maxOrderNumber(P^*)$ 
166    $P' := \{\langle \text{PREPARE}, r, v', o, m \rangle_{\tau(r, \mathfrak{D}, v|o, -)} \mid o > o\_base_r \wedge \langle \text{PREPARE}, \cdot, \cdot, o, m \rangle \dots \in P^*\}$ 
167   for  $o : o\_base_r < o \leq o\_max$  do
168      $pr' := \langle \text{PREPARE}, r, v', o, m \rangle \dots \mid pr' \in P'$ 
169     invoke  $tss_r.createIndependentCertificate(\mathfrak{D}, v'|o, pr')$ 
170   return  $P'$ 

172 upon global call  $isCorrectViewTransition(P', v' \in {}^v\mathbb{N}, o\_base \in {}^o\mathbb{N}, B \subseteq AbortState)$  do
173    $P^* := \{pr \mid \langle P, \cdot \rangle \in B \wedge pr \in P\}$ 
174   return  $isCorrectPrepareSet(P', v', o\_base, \cdot)$ 
175      $\wedge (\forall pr \in P^* : pr = \langle \text{PREPARE}, \cdot, \cdot, o, m \rangle \dots \wedge (o > o\_base \rightarrow \langle \text{PREPARE}, \cdot, v, o, m \rangle \dots \in P'))$ 
176      $\wedge (\forall pr' \in P' : pr' = \langle \text{PREPARE}, \cdot, \cdot, o, m \rangle \dots \wedge \langle \text{PREPARE}, \cdot, \cdot, o, m \rangle \dots \in P^*)$ 

178 upon input call  $enterView(v' \in {}^v\mathbb{N}, P' \subseteq Message)$  do
179    $ostatus_r.enterView(v')$ 
180   if  $|P'| = 0 \vee r \neq ostatus_r.lastestLeader()$  then
181      $o\_prep_r := o\_base_r$ 
182   else
183      $o\_max := \maxOrderNumber(P)$ 
184      $o\_prep_r := o\_max$ 
185      $o\_act_r := o\_max$ 
186      $o\_comm_r := o\_base_r$ 
187   output  $repgrpconn_r.remove(\{m \mid m = \langle \text{PREPARE}, \dots \rangle \dots \vee m = \langle \text{COMMIT}, \dots \rangle \dots\})$ 
188    $installOrderingState(P')$ 

190 upon internal call  $orderNumber(m \in Message)$ 
191   asserts
192      $m = (\langle \text{PREPARE}, \cdot, \cdot, o, \cdot \rangle \dots \vee \langle \text{COMMIT}, \cdot, \cdot, o, \cdot \rangle \dots)$ 
193   do
194     return  $o$ 

196 upon internal call  $\maxOrderNumber(I \subseteq Message)$ 
197   return  $\max(\{o \mid m \in I \wedge o = orderNumber(m)\})$ 

199 upon internal call  $\minOrderNumber(I \subseteq Message)$ 
200   return  $\min(\{o \mid m \in I \wedge o = orderNumber(m)\})$ 

```


4.7 Execution

Listing 17: Specification of the service execution.

```

1 module Replicar::Execution

3 uses cmr ∈ Replicar::Client
4 uses kmr ∈ Replicar::Checkpointing

6 def global ServiceSnapshot  $\stackrel{\text{def}}{=} \mathcal{S} \times (\mathcal{C} \rightarrow \text{ReturnEntry})$ 

8 upon init() do
9   svcr := Service()
10  oexecr ∈ oℕ := 0

12 upon input call executeRequest(o ∈ oℕ, rq = ⟨REQUEST, c, t, svccmd⟩...)
13   asserts
14     o = oexecr + 1
15   do
16     if ¬cmr.isRequestExecuted(rq) then
17       svcret := invoke svcr.invoke(c, svccmd)
18       output cmr.requestExecuted(rq, svcret)
19       oexecr := o
20       output kmr.stateReached(o)

22 upon output call createSnapshot() do
23   s := input svcr.createStateSnapshot()
24   retvals := input cmr.createReturnValueMapSnapshot()
25   return ⟨s, retvals⟩

27 upon input call installSnapshot(oexec ∈ oℕ, ⟨s, retvals⟩ ∈ ServiceSnapshot) do
28   invoke svcr.installStateSnapshot(s)
29   invoke cmr.installReturnValueMap(retvals)
30   oexecr := oexec

```

4.8 Checkpointing Protocol

Listing 18: Specification of the checkpointing protocol.

```

1 module Replicar::Checkpointing

3 uses omr ∈ Replicar::Ordering
4 uses emr ∈ Replicar::Execution

6 def global
7   CheckpointState  $\stackrel{\text{def}}{=} {}^o\mathbb{N} \times \text{ServiceSnapshot}$ 
8   CheckpointSet  $\stackrel{\text{def}}{=} \{X \mid \exists o \forall ck \in X: ck = \langle \text{CHECKPOINT}, \cdot, o', \cdot \rangle \dots \wedge o = o'\}$ 
9   ProvenCheckpoint  $\stackrel{\text{def}}{=} \text{CheckpointSet} \times \text{ServiceSnapshot}$ 

11 let
12   chkpt_int ∈ ℕ1, chkpt_int ≤ orderWindowSize()
13   retent_wnd ∈ ℕ, retent_wnd < orderWindowSize()

15 upon global call checkpointBase(o)
16   return if o < retent_wnd then 0 else o − retent_wnd

18 upon global call checkpointNumber(⟨o, s*⟩ ∈ CheckpointState) do
19   return o

21 upon global call checkpointState(⟨o, s*⟩ ∈ CheckpointState) do
22   return s*

24 upon global call checkpointNumber(K ∈ CheckpointSet)
25   asserts
26     ∃o ∀⟨CHECKPOINT, ·, o', ·⟩... ∈ K: o = o'
27   do
28     return o

30 upon init() do
31   kmsgsr ⊆ Message := {}
32   chkptsr ⊆ CheckpointState := {}
33   o_stabr ∈ oℕ := 0

35 upon input call stateReached(o ∈ oℕ) do
36   if (o mod chkpt_int) = 0 then
37     createCheckpoint(o)

39 upon internal call createCheckpoint(o ∈ oℕ) do
40   ks := ⟨o, emr.createSnapshot()⟩
41   chkptsr := chkptsr ∪ {ks}
42   sendCheckpoint(ks)

44 upon internal call sendCheckpoint(⟨o, s*⟩ ∈ CheckpointState)
45   ck := ⟨CHECKPOINT, r, o, digest(s*)⟩τ(r, ℳ, 0, 0)
46   invoke tssr.createConsecutiveCertificate(ℳ, 0, ck)
47   kmsgsr := kmsgsr ∪ {ck}
48   output repgrpconnr.send(ck)

50 upon global call isCorrectCheckpoint(ck, r' ∈ ℛ, o ∈ oℕ, d ∈ ⟨DIGEST...⟩) do
51   return ck = ⟨CHECKPOINT, r', o, d⟩τ(r', ℳ, 0, 0) ∧ hasValidCertificate(ck)

53 upon input call receive(ck' = ⟨CHECKPOINT, r', o, d⟩...) do

```

```

54  if isRelevantCheckpoint(ck')  $\wedge$  isCorrectCheckpoint(ck', ...)
     $\wedge$  ( $\nexists d' : \langle \text{CHECKPOINT}, r', o, d' \rangle \dots \in kmsgs_r \wedge d' \neq d$ ) then
55      storeCheckpoint(ck')

57 upon internal call isRelevantCheckpoint( $\langle \text{CHECKPOINT}, r', o, \cdot \rangle \dots$ ) do
58   return  $r' \neq r \wedge isRelevantCheckpoint(o)$ 

60 upon internal call isRelevantCheckpoint( $o \in {}^o\mathbb{N}$ ) do
61   return  $o = \text{input } om_r.orderWindowBase() \vee \text{input } om_r.inOrderWindow(o)$ 

63 upon internal call storeCheckpoint( $ck' = \langle \text{CHECKPOINT}, r', o, d \rangle \dots$ ) do
64    $kmsgs_r := kmsgs_r \setminus \{ck \in kmsgs_r \mid ck = \langle \text{CHECKPOINT}, r', o, \cdot \rangle \dots\}$ 
65    $kmsgs_r := kmsgs_r \cup \{ck'\}$ 

67 upon internal call discardCheckpoints() do
68    $o\_base := \text{input } om_r.orderWindowBase()$ 
69    $kmsgs_r := kmsgs_r \setminus \{\langle \text{CHECKPOINT}, \cdot, o, \cdot \rangle \dots \in kmsgs_r \mid o < o\_base\}$ 
70   output regrprconnr.remove( $\{ck \mid ck = \langle \text{CHECKPOINT}, \dots \rangle \dots \wedge ck \notin kmsgs_r\}$ )

72 upon global call isCorrectCheckpointCertificate( $K \subseteq Message, o \in {}^o\mathbb{N}, d \in \langle \text{DIGEST} \dots \rangle$ ) do
73   return  $o = 0 \wedge |K| = 0 \vee |K| > f \wedge (\forall ck \in K : isCorrectCheckpoint(ck, \cdot, o, d))$ 

75 upon internal task checkpointStable( $o' \in {}^o\mathbb{N}, d \in \langle \text{DIGEST} \dots \rangle$ )
76   with
77      $o' > o\_stab_r \wedge (\exists s^* : digest(s^*) = d \wedge \langle o', s^* \rangle \in chkpts_r)$ 
78      $\exists K' \subseteq kmsgs_r : isCorrectCheckpointCertificate(K', o', d)$ 
79      $\wedge (\forall K \subseteq kmsgs_r \exists o : isCorrectCheckpointCertificate(K, o, \cdot) \rightarrow o \leq o')$ 
80   do
81     forwardWindow( $o'$ )

83 upon internal call forwardWindow( $o' \in {}^o\mathbb{N}$ ) do
84    $o\_stab_r := o'$ 
85    $chkpts_r := chkpts_r \setminus \{\langle o, \cdot \rangle \in chkpts_r \mid o < o'\}$ 
86   if checkpointBase( $o'$ )  $> \text{input } om_r.orderWindowBase()$  then
87     invoke omr.forwardOrderWindow(checkpointBase( $o'$ ))
88     discardCheckpoints()

90 upon global call isCorrectProvenCheckpoint( $pk, o \in {}^o\mathbb{N}, s^* \in ServiceSnapshot$ ) do
91   return  $pk = \langle C, s^* \rangle \in ProvenCheckpoint \wedge isCorrectCheckpointCertificate(C, o, digest(s^*))$ 

93 upon output call stableCheckpointNumber() do
94   return  $o\_stab_r$ 

96 upon internal call hasStableCheckpoint() do
97   return  $o\_stab_r > 0$ 

99 upon internal call stableCheckpoint()  $\rightarrow ProvenCheckpoint$ 
100  asserts
101    hasStableCheckpoint()
102  do
103     $\langle o, s^* \rangle \in chkpts_r \mid o = o\_stab_r$ 
104     $K := \{\langle \text{CHECKPOINT}, \cdot, o, digest(s^*) \rangle \dots \in kmsgs_r\}$ 
105    return  $\langle K, s^* \rangle$ 

107 upon output call stableCheckpointCertificate()  $\rightarrow CheckpointSet$  do
108   if  $\neg hasStableCheckpoint()$  then
109     return  $\emptyset$ 

```

```

110 else
111    $\langle K, \cdot \rangle := \text{stableCheckpoint}()$ 
112   return  $K$ 

114 upon output call  $\text{obtainCheckpointingState}(o\_stab \in {}^o\mathbb{N}, o\_comm \in {}^o\mathbb{N})$  do
115   if  $o\_comm < \text{input } om_r.\text{orderWindowBase}()$  then
116     return  $\text{stableCheckpoint}()$ 
117   else if  $o\_stab < o\_stab_r$  then
118     return  $\text{stableCheckpointCertificate}()$ 
119   else
120     return  $\emptyset$ 

122 upon global call  $\text{isCorrectCheckpointingState}(ks, o\_stab \in {}^o\mathbb{N})$  do
123   return  $ks = \emptyset$ 
124    $\vee ks \in \text{CheckpointSet} \wedge (|ks| = 0 \vee \text{isCorrectCheckpointCertificate}(ks, o\_stab, \cdot))$ 
125    $\vee \text{isCorrectProvenCheckpoint}(ks, o\_stab, \cdot)$ 

127 upon input call  $\text{installCheckpointingState}(ks)$ 
128   asserts
129      $\text{isCorrectCheckpointingState}(ks, \dots)$ 
130   do
131     if  $ks \in \text{CheckpointSet}$  then
132        $\text{installCheckpointCertificate}(ks)$ 
133     else if  $ks = \langle K, s^* \rangle \wedge \neg \text{input } om_r.\text{isCommitted}(\text{checkpointNumber}(K))$  then
134        $o := \text{checkpointNumber}(K)$ 
135       output  $em_r.\text{installSnapshot}(o, s^*)$ 
136        $\text{forwardWindow}(o)$ 
137        $\text{createCheckpoint}(o)$ 
138        $\text{installCheckpointCertificate}(K)$ 

140 upon input call  $\text{installCheckpointCertificate}(K \in \text{CheckpointSet})$ 
141   asserts
142      $\exists o: \text{isCorrectCheckpointCertificate}(K, o, \cdot)$ 
143   do
144     if  $\text{isRelevantCheckpoint}(o)$  then
145       for all  $ck \in K$  do
146         if  $\text{isRelevantCheckpoint}(ck)$  then
147            $\text{storeCheckpoint}(ck)$ 
148       if  $o > o\_stab_r$  then
149         check  $\text{checkpointStable}()$ 

```

4.9 View-Change Protocol

Listing 19: Specification of the view-change protocol.

```

1 module Replicar::ViewChange

3 uses cmr ∈ Replicar::Client
4 uses omr ∈ Replicar::Ordering
5 uses kmr ∈ Replicar::Checkpointing
6 uses smr ∈ Replicar::StateTransfer

8 upon init() do
9   vstatusr := ViewStatus(r)
10  v_stabr  $\stackrel{\text{def}}{=}$  vstatusr.stableView()
11  v_curr  $\stackrel{\text{def}}{=}$  vstatusr.currentView()
12  vmsgsr ⊆ Message := {}
13  nv_timeoutr := AdaptiveTimeout( $\mathcal{R} \setminus \{r\}$ )

15 upon internal call stableViewCertificate() do
16   asserts
17   v_stabr > 0
18   do
19     return ⟨NEW-VIEW, ·, v_stabr, ...⟩... ∈ vmsgsr

21 upon internal call suspectsAcceptedLeader() do
22   return cmr.suspectsLeader()

24 upon internal call isLeaderSuspectedByQuorum() do
25   return ∃V ⊆ vmsgsr: |V| > f
26   ∧ (∀vc ∈ V ∃v_to: v_to > v_curr ∧ isCorrectViewChange(vc, ·, ·, v_to))

28 upon internal task acceptedLeaderSuspected()
29   with
30     vstatusr.isInStableView()
31     ∧ (suspectsAcceptedLeader() ∨ isLeaderSuspectedByQuorum())
31   do
32     leaveViewFor(v_curr + 1)

34 upon internal call leaveViewFor(v' ∈  ${}^v\mathbb{N}$ ) do
35   v := v_curr
36   l := leader(v)
37   K := input kmr.stableCheckpointCertificate()
38   b := input omr.createAbortState()
39   vc := ⟨VIEW-CHANGE, r, v_stabr, v', K, b⟩...
40   if nv_timeoutr.isRunning() then
41     nv_timeoutr.stop()
42   output omr.abortView(v', vc)
43   vstatusr.leaveViewFor(v')
44   storeViewChange(vc)
45   output repgrpconnr.send(vc)
46   if v = v_stabr then
47     output smr.startStateRequestTask()

49 upon global call isCorrectViewChange(vc, r' ∈  $\mathcal{R}$ , v_from ∈  ${}^v\mathbb{N}$ , v_to ∈  ${}^v\mathbb{N}$ ) do
50   return ∃K ⊆ Message ∃b ∈ AbortState ∃o_base, o_stab ∈  ${}^o\mathbb{N}$ :
51     vc = ⟨VIEW-CHANGE, r', v_from, v_to, K, b⟩... ∧ v_from < v_to
52     ∧ isCorrectCheckpointCertificate(K, o_stab, ·)

```

```

53       $\wedge o\_base = \text{checkpointBase}(o\_stab)$ 
54       $\wedge \text{isCorrectAbortState}(b, v\_from, v\_to, o\_base, vc)$ 

56 upon input call  $\text{receive}(vc' = \langle \text{VIEW-CHANGE}, r', v\_from, v\_to, K, b \rangle \dots)$  do
57   if  $\text{isRelevantViewChange}(vc') \wedge \text{isCorrectViewChange}(vc', \dots)$  then
58      $\text{installViewChange}(vc')$ 

60 upon internal call  $\text{isRelevantViewChange}(vc' = \langle \text{VIEW-CHANGE}, r', v\_from, v\_to, \dots \rangle \dots)$  do
61   return  $r' \neq r \wedge (v\_to > v\_cur_r \vee \neg vstatus_r.\text{isInStableView} \wedge v\_to = v\_cur_r)$ 

63 upon internal call  $\text{storeViewChange}(vc' = \langle \text{VIEW-CHANGE}, \cdot, \cdot, v\_to', \dots \rangle \dots)$  do
64    $vmsgs_r := vmsgs_r \cup \{vc'\}$ 
65    $\text{discardViewChanges}()$ 

67 upon internal call  $\text{discardViewChanges}()$  do
68   if  $vstatus_r.\text{isInStableView}()$  then
69      $VC^{next} := \emptyset$ 
70   else
71      $VC^{next} := \{vc \in vmsgs_r \mid vc = \langle \text{VIEW-CHANGE}, \cdot, \cdot, v\_cur_r, \dots \rangle \dots\}$ 
72      $VC^{max} := \{vc\_max \in vmsgs_r \mid vc\_max = \langle \text{VIEW-CHANGE}, r', \cdot, v\_max, \dots \rangle \dots$ 
73  $\wedge v\_max > v\_cur_r$ 
74  $\wedge (\forall vc \in vmsgs_r: vc = \langle \text{VIEW-CHANGE}, r', \cdot, v, \dots \rangle \dots \rightarrow v \leq v\_max)\}$ 
75      $VS := \{\langle v, V^* \rangle \mid V^* \subseteq vmsgs_r \wedge v > v\_stab_r$ 
76  $\wedge \text{isCorrectViewChangeCertificate}(V^*, v\_stab_r, v)$ 
77  $\wedge (\forall V \subseteq vmsgs_r:$ 
78  $\text{isCorrectViewChangeCertificate}(V, v\_stab_r, v) \rightarrow |V| \leq |V^*|)\}$ 
79     if  $v\_cur_r \leq v\_stab_r + 1$  then
80        $V^{cur} := \emptyset$ 
81     else if  $\langle v\_cur_r, V^* \rangle \in VS$  then
82        $V^{cur} := V^*$ 
83     else
84        $V^{cur} := V^* \mid \langle v\_cur_r - 1, V^* \rangle \in VS$ 
85     if  $|VS| = 0$  then
86        $V^{max} := \emptyset$ 
87     else
88        $V^{max} := V^* \mid \langle v\_max, V^* \rangle \in VS \wedge (\forall \langle v, \cdot \rangle \in VS: v \leq v\_max)$ 
89      $vmsgs_r := vmsgs_r \setminus \{vc \in vmsgs_r \mid vc = \langle \text{VIEW-CHANGE}, \dots \rangle \dots$ 
90  $\wedge vc \notin (VC^{next} \cup VC^{max} \cup V^{cur} \cup V^{max})\}$ 
91     output  $\text{repgprconn}_r.\text{remove}(\{vc \mid vc = \langle \text{VIEW-CHANGE}, \dots \rangle \dots$ 
92  $\wedge vc \notin (VC^{next} \cup VC^{max} \cup V^{cur} \cup V^{max})\})$ 

94 upon input call  $\text{installViewChange}(vc' = \langle \text{VIEW-CHANGE}, \cdot, v\_from, v\_to, K, b \rangle \dots)$  do
95    $\text{storeViewChange}(vc')$ 
96    $\text{abortStateReceived}(v\_from, K, b)$ 

98 upon internal call  $\text{abortStateReceived}(v \in {}^v\mathbb{N}, K \in \text{CheckpointSet}, b \in \text{AbortState})$  do
99   output  $km_r.\text{installCheckpointCertificate}(K)$ 
100  if  $v = v\_stab_r$  then
101    output  $om_r.\text{installAbortStaet}(b)$ 

103 upon global call  $\text{isCorrectViewChangeCertificate}(V \subseteq \text{Message}, v\_from \in {}^v\mathbb{N}, v\_to \in {}^v\mathbb{N})$  do
104  return  $|V| \geq q \wedge (\forall vc \in V \exists v\_from^{\leq}:$ 
 $v\_from^{\leq} \leq v\_from \wedge \text{isCorrectViewCange}(vc, \cdot, v\_from^{\leq}, v\_to)$ 

106 upon internal call  $\text{isInOrderWindow}(K \subseteq \text{Message})$ 
107  return  $\text{checkpointBase}(\text{checkpointNumber}(K)) \leq \text{input } om_r.\text{orderWindowBase}()$ 

```

```

109 upon internal call isInOrderWindow( $V \subseteq \text{Message}$ ,  $A \subseteq \text{Message}$ )
110 return checkpointBase(maxCheckpointNumber( $V, A$ ))  $\leq$  input omr.orderWindowBase()

112 upon internal call maxCheckpointNumber( $V \subseteq \text{Message}$ ,  $A \subseteq \text{Message}$ ) do
113 return max( $\{ \text{checkpointNumber}(K) \mid$ 
     $\langle \text{VIEW-CHANGE}, \cdot, \cdot, \cdot, K, \dots \rangle \dots \in V \vee \langle \text{NEW-VIEW-ACK}, \cdot, \cdot, K, \cdot \rangle \dots \}$ )

115 upon internal task newViewExpected()
116 with
117    $\neg \text{vstatus}_r.\text{isInStableView}() \wedge l' = \text{leader}(v\_cur_r) \wedge l' \neq r$ 
118    $\exists V \subseteq \text{msgs}_r: \text{isCorrectViewChangeCertificate}(V, v\_stab_r, v\_cur_r)$ 
119    $\neg \text{nv\_timeout}_r.\text{isRunning}()$ 
120 do
121   nv\_timeoutr.start( $l'$ )

123 upon internal call suspectsDesignatedLeader() do
124 return nv\_timeoutr.isExpired()

126 upon internal task designatedLeaderSuspected()
127 with
128    $\neg \text{vstatus}_r.\text{isInStableView}() \wedge (\text{suspectsDesignatedLeader}() \vee \text{isLeaderSuspectedByQuorum}())$ 
129    $\exists V \subseteq \text{msgs}_r: \text{isCorrectViewChangeCertificate}(V, v\_stab_r, v\_cur_r) \wedge \text{isInOrderWindow}(V, \emptyset)$ 
130 do
131   installAbortStates( $V$ )
132   leaveViewFor( $v\_cur_r + 1$ )

134 upon internal task viewMissed( $v' \in {}^v\mathbb{N}$ )
135 with
136    $\exists V \subseteq \text{msgs}_r \exists v: v' > v\_cur_r \wedge \text{isCorrectViewChangeCertificate}(V, v\_stab_r, v')$ 
     $\wedge \text{isInOrderWindow}(V, \emptyset)$ 
137 do
138   installAbortStates( $V$ )
139   leaveViewFor( $v'$ )

141 upon internal call installAbortStates( $V \subseteq \text{Message}$ )
142 for all  $vc \in V$  do
143    $vc := \langle \text{VIEW-CHANGE}, \cdot, v\_from, \dots, b \rangle \dots$ 
144   if  $v\_from = v\_stab_r$  then
145     output omr.installAbortState( $b$ )

147 upon global call isCorrectNewViewCertificate( $V \subseteq \text{Message}$ ,  $A \subseteq \text{Message}$ ,  $v\_from \in {}^v\mathbb{N}$ ,  $v\_to \in {}^v\mathbb{N}$ )
    do
148 return isCorrectViewChangeCertificate( $V, v\_from, v\_to$ )
     $\wedge (\forall na \in A: \text{isCorrectNewViewAck}(na, \cdot, v\_from))$ 
149    $\wedge |\{r' \mid \langle \text{VIEW-CHANGE}, r', v\_from, \dots \rangle \dots \in V \vee \langle \text{NEW-VIEW-ACK}, r', \dots \rangle \dots \in A\}| > f$ 

152 upon internal task newLeaderViewReady()
153 with
154    $\exists v' \geq v\_cur_r \exists V, A \subseteq \text{msgs}_r: r = \text{leader}(v')$ 
155    $\wedge \text{isCorrectNewViewCertificate}(V, A, \cdot, v')$ 
156    $\wedge \text{isInOrderWindow}(V, A)$ 
157 do
158   establishView( $v', V, A$ )

160 upon internal call establishView( $v' \in {}^v\mathbb{N}$ ,  $V \subseteq \text{Message}$ ,  $A \subseteq \text{Message}$ ) do
161 if  $v' > v\_cur_r$  then
162   output omr.abortView( $v'$ )

```

```

163    $v_{cur_r} := v'$ 
164    $K := \text{input } km_r.\text{stableCheckpointCertificate}()$ 
165    $B := \text{allAbortStates}(V, A)$ 
166    $P' := \text{input } om_r.\text{createViewTransition}(v', B)$ 
167    $nv := \langle \text{NEW-VIEW}, r, v', K, V, A, P' \rangle_{\tau(r, \mathfrak{N}, v', -)}$ 
168   invoke  $tss_r.\text{createIndependentCertificate}(\mathfrak{N}, v', nv)$ 
169    $\text{storeNewView}(nv)$ 
170    $\text{enterView}(v', P')$ 
171   output  $\text{repgrpconn}_r.\text{send}(nv)$ 

173 upon global call  $\text{isCorrectNewView}(nv, v \in {}^v\mathbb{N}, K \subseteq \text{Message})$  do
174   return  $\exists V, A, P', o\_stab, o\_base:$ 
175      $nv = \langle \text{NEW-VIEW}, \text{leader}(v), v, K, V, P' \rangle_{\tau(l, \mathfrak{N}, v, -)}$ 
176      $\wedge \text{hasValidCertificate}(nv)$ 
177      $\wedge \text{isCorrectNewViewCertificate}(V, A, \cdot, v)$ 
178      $\wedge \text{isCorrectCheckpointCertificate}(K, o\_stab, \cdot)$ 
179      $\wedge o\_stab \geq \text{maxCheckpointNumber}(V, A)$ 
180      $\wedge o\_base = \text{checkpointBase}(o\_stab)$ 
181      $\wedge \text{isCorrectViewTransition}(P', v, o\_base, \text{allAbortStates}(V, A))$ 

183 upon input call  $\text{receive}(nv' = \langle \text{NEW-VIEW}, \dots \rangle \dots)$  do
184   if  $\text{isRelevantNewView}(nv') \wedge \text{isCorrectNewView}(nv', \dots)$  then
185      $\text{installNewView}(nv')$ 

187 upon internal call  $\text{isRelevantNewView}(nv' = \langle \text{NEW-VIEW}, \cdot, v', \dots \rangle \dots)$  do
188   return  $v' > \text{latestKnownView}()$ 

190 upon internal call  $\text{storeNewView}(nv' = \langle \text{NEW-VIEW}, \dots \rangle \dots)$  do
191    $\text{msgs}_r := \text{msgs}_r \cup \{nv'\}$ 
192    $\text{discardNewViews}()$ 

194 upon internal call  $\text{discardNewViews}()$  do
195    $NV^{cur} := \{nv\_cur \in \text{msgs}_r \mid nv\_cur = \langle \text{NEW-VIEW}, \cdot, v\_cur, \dots \rangle \dots$ 
196      $\wedge (v\_cur = v\_stab_r \vee v\_cur = \text{latestKnownView}())\}$ 
197    $\text{msgs}_r := \text{msgs}_r \setminus \{nv \in \text{msgs}_r \mid nv = \langle \text{NEW-VIEW}, \dots \rangle \dots \wedge nv \notin NV^{cur}\}$ 
198   output  $\text{repgrpconn}_r.\text{remove}(\{nv \mid nv = \langle \text{NEW-VIEW}, \dots \rangle \dots \wedge nv \notin NV^{cur}\})$ 

199 upon internal call  $\text{installNewView}(nv' = \langle \text{NEW-VIEW}, l, v', K, V, P' \rangle \dots)$  do
200    $\text{storeNewView}(nv')$ 
201   for all  $vc \in V$  do
202     if  $\text{isRelevantViewChange}(vc)$  then
203        $\text{installViewChange}(vc)$ 
204   output  $km_r.\text{installCheckpointCertificate}(K)$ 
205   for all  $\langle \text{VIEW-CHANGE}, \dots, K', \dots \rangle \dots \in V$ 
206     output  $km_r.\text{installCheckpointCertificate}(K')$ 

208 upon internal call  $\text{allAbortStates}(V \subseteq \text{Message}, A \subseteq \text{Message})$  do
209   return  $\{b \mid \langle \text{VIEW-CHANGE}, \dots, b \rangle \dots \in V \vee \langle \text{NEW-VIEW-ACK}, \dots, b \rangle \dots \in A\}$ 

211 upon internal call  $\text{latestKnownView}()$  do
212   return  $\max(\{v \mid \langle \text{NEW-VIEW}, \cdot, v, \dots \rangle \dots \in \text{msgs}_r\} \cup \{0\})$ 

214 upon internal task  $\text{newFollowerViewReady}(nv = \langle \text{NEW-VIEW}, l, v', K, \cdot, \cdot, P' \rangle \dots)$ 
215   with
216      $nv \in \text{msgs}_r \wedge v' = \text{latestKnownView}() \wedge v' > v\_stab_r \wedge \text{isInOrderWindow}(K)$ 
217   do
218      $\text{isbelated} := v' < v\_cur_r$ 

```



```

219   if  $v' > v\_cur_r$  then
220     output  $om_r.abortView(v')$ 
221      $v\_cur_r := v'$ 
222      $enterView(v', P')$ 
223     if  $isbelated$  then
224        $sendNewViewAck()$ 

226 upon internal call  $sendNewViewAck()$  do
227    $K := \text{input } km_r.stableCheckpointCertificate()$ 
228    $b := \text{input } om_r.createAbortState(v\_cur_r)$ 
229    $na := \langle \text{NEW-VIEW-ACK}, r, v\_cur_r, K, b \rangle_{\tau(r, \mathfrak{M}, 0, 0)}$ 
230   invoke  $tss_r.createContinuingCertificate(\mathfrak{M}, 0, na)$ 
231    $storeNewViewAck(na)$ 
232   output  $repgrpconn_r.send(na)$ 

234 upon global call  $isCorrectNewViewAck(na, r' \in \mathcal{R}, v \in {}^v\mathbb{N})$  do
235   return  $\exists K \subseteq Message \exists b \in AbortState \exists o\_base, o\_stab \in {}^o\mathbb{N}:$ 
236      $na = \langle \text{NEW-VIEW-ACK}, r', v, K, b \rangle_{\tau(r, \mathfrak{M}, 0, 0)}$ 
237      $\wedge hasValidCertificate(na)$ 
238      $\wedge isCorrectCheckpointCertificate(K, o\_stab, \cdot)$ 
239      $\wedge o\_base = checkpointBase(o\_stab)$ 
240      $\wedge isCorrectAbortState(b, v\_from, v\_to, o\_base)$ 

242 upon input call  $receive(na' = \langle \text{NEW-VIEW-ACK}, \dots \rangle \dots)$  do
243   if  $isRelevantNewViewAck(na') \wedge isCorrectNewViewAck(na', \dots)$  then
244      $installNewViewAck(na')$ 

246 upon internal call  $isRelevantNewViewAck(na' = \langle \text{NEW-VIEW-ACK}, r', v', \dots \rangle \dots)$  do
247   return  $v' \geq v\_stab_r \wedge (\forall na \in vmsgs_r: na = \langle \text{NEW-VIEW-ACK}, r', v, \dots \rangle \dots \rightarrow v' > v)$ 

249 upon internal call  $storeNewViewAck(na' = \langle \text{NEW-VIEW-ACK}, \dots \rangle \dots)$  do
250    $vmsgs_r := vmsgs_r \cup \{na\}$ 
251    $discardNewViewAcks()$ 

253 upon internal call  $discardNewViewAcks()$  do
254    $NA^{max} := \{na\_max \in vmsgs_r \mid na\_max = \langle \text{NEW-VIEW-ACK}, r', v\_max, \dots \rangle \dots$ 
255    $\wedge v\_max \geq v\_stab_r$ 
256    $\wedge (\forall na \in vmsgs_r: na = \langle \text{NEW-VIEW-ACK}, r', v, \dots \rangle \dots \rightarrow v \leq v\_max)\}$ 
257    $vmsgs_r := vmsgs_r \setminus \{na \in vmsgs_r \mid na = \langle \text{NEW-VIEW-ACK}, \dots \rangle \dots \wedge na \notin NA^{max}\}$ 
258   output  $repgrpconn_r.remove(\{na \mid na = \langle \text{NEW-VIEW-ACK}, \dots \rangle \dots \wedge na \notin NA^{max}\})$ 

260 upon internal call  $installNewViewAck(na' = \langle \text{NEW-VIEW-ACK}, \cdot, v', K, b \rangle \dots)$  do
261    $storeNewViewAck(na')$ 
262    $abortStateReceived(v', K, b)$ 

264 upon internal call  $enterView(v' \in {}^v\mathbb{N}, P' \subseteq Message)$  do
265    $vstatus_r.enterView(v')$ 
266    $discardViewChanges()$ 
267    $discardNewViews()$ 
268    $discardNewViewAcks()$ 
269   if  $vstatus_r.isInStableView()$  then
270     if  $nv\_timeout_r.isRunning()$  then
271        $nv\_timeout_r.stop()$ 
272     output  $sm_r.stopStateRequestTask()$ 
273   output  $om_r.enterView(v', P')$ 
274   output  $cm_r.enterView(v')$ 

```

```

276 upon output call createStateRequestContent() do
277   o_comm := input omr.lastCommitted()
278   o_stab := input kmr.stableCheckpointNumber()
279   return  $\langle \text{latestKnownView}(), v_{\text{cur}}, o_{\text{stab}}, o_{\text{comm}} \rangle$ 

281 upon global call isCorrectStateRequestContent(sr) do
282   return  $sr = \langle v_{\text{latest}} \in {}^v\mathbb{N}, v_{\text{cur}} \in {}^v\mathbb{N}, o_{\text{stab}} \in {}^o\mathbb{N}, o_{\text{comm}} \in {}^o\mathbb{N} \rangle \wedge o_{\text{stab}} \leq o_{\text{comm}}$ 

284 upon output call obtainState( $\langle v_{\text{latest}}, v_{\text{cur}}, o_{\text{stab}}, o_{\text{comm}} \rangle$ ) do
285   if  $v_{\text{latest}} > v_{\text{stab}_r} \vee v_{\text{stab}_r} \neq \text{latestKnownView}()$ 
286     return  $\emptyset$ 
287   if  $v_{\text{latest}} < v_{\text{stab}_r}$  then
288     nv := stableViewCertificate()
289     I := input omr.obtainOrderingState( input omr.orderWindowBase())
290   else
291     nv :=  $\emptyset$ 
292     I := input omr.obtainOrderingState(o_comm)
293   V :=  $\{vc \in \text{msgs}_r \mid vc = \langle \text{VIEW-CHANGE}, \dots, v_{\text{to}}, \cdot \rangle \dots$ 
      $\wedge (v_{\text{to}} > v_{\text{latest}} \geq v_{\text{cur}} \vee v_{\text{to}} \geq v_{\text{cur}} > v_{\text{latest}}) \}$ 
294   A :=  $\{na \in \text{msgs}_r \mid na = \langle \text{NEW-VIEW-ACK}, \cdot, v, \dots \rangle \dots \wedge v \geq v_{\text{latest}} \}$ 
295   ks := input kmr.obtainCheckpointingState(o_stab, o_comm),
296   return if  $nv \neq \emptyset \vee |V| > 0 \vee |A| > 0 \vee ks \neq \emptyset \vee |I| > 0$ 
297     then  $\langle v_{\text{stab}_r}, nv, V, A, ks, I \rangle$  else  $\emptyset$ 

299 upon global call isCorrectState(state)
300   return  $state = \langle v, nv, V, A, ks, I \rangle \wedge (\exists K', o_{\text{base}}, o_{\text{stab}}:$ 
301      $\wedge (nv = \emptyset \vee \text{isCorrectNewView}(nv, v, K', \dots) \wedge o_{\text{stab}} \geq \text{checkpointNumber}(K'))$ 
302      $\wedge (\forall vc \in V \exists v_{\text{to}}: \text{isCorrectViewChange}(vc, \cdot, \cdot, v_{\text{to}}) \wedge v_{\text{to}} > v)$ 
303      $\wedge (\forall na \in A \exists v': \text{isCorrectNewViewAck}(na, \cdot, v') \wedge v' \geq v)$ 
304      $\wedge o_{\text{base}} = \text{checkpointBase}(o_{\text{stab}})$ 
305      $\wedge \text{isCorrectCheckpointingState}(ks, o_{\text{stab}})$ 
306      $\wedge \text{isCorrectOrderingState}(I, v, o_{\text{base}}))$ 

308 upon input call installState( $\langle v, nv, V, A, ks, I \rangle$ )
309   output kmr.installCheckpointingState(ks)
310   if isRelevantNewView(nv) then
311     installNewView(nv)
312   check newFollowerViewReady()
313   if  $v = v_{\text{stab}_r}$  then
314     output omr.installOrderingState(I)
315   for all  $vc \in V$  do
316     if isRelevantViewChange(vc) then
317       installViewChange(vc)
318   for all  $na \in A$  do
319     if isRelevantNewViewAck(na) then
320       installNewViewAck(na)

```

4.10 State Transfer Protocol

Listing 20: Specification of the state-transfer protocol.

```

1 module Replicar::StateTransfer

3 uses vmr ∈ Replicar::ViewChange

5 let state_periodr ∈ ℕ

7 upon init() do
8   state_timerr := Timer()

12 upon input call startStateRequestTask() do
13   output state_timerr.schedule(state_periodr)

14 upon input call stopStateRequestTask() do
15   output state_timerr.cancel()

16 upon internal task periodicallyRequestState()
17   with
18     input state_timerr.isExpired()
19   do
20     requestState()
21     state_timerr.schedule(state_periodr)

22 upon input call requestState() do
23   src := input vmr.createStateRequestContent()
24   sr := ⟨STATE-REQUEST, r, src⟩τ(r, ℳ, 0, 0)
25   invoke tssr.createConsecutiveCertificate(ℳ, 0, sr)
26   output repgrpconnr.remove({m | m = ⟨STATE-REQUEST, ...⟩...})
27   output repgrpconnr.send(sr)

30 upon global call isCorrectStateRequest(sr, r' ∈ ℛ, src) do
31   return sr = ⟨STATE-REQUEST, r', src⟩τ(r', ℳ, 0, 0)
      ∧ hasValidCertificate(sr) ∧ isCorrectStateRequestContent(src)

32 upon input call receive(rs = ⟨STATE-REQUEST, r', src⟩...) do
33   if isCorrectStateRequest(rs, r', src) then
34     state := input vmr.obtainState(src)
35     if state ≠ ∅
36       sendState(r', state)

37 upon internal call sendState(r' ∈ ℛ, state) do
38   sm := ⟨STATE, r, state⟩τ(r, ℳ, 0, 0)
39   invoke tssr.createConsecutiveCertificate(ℳ, 0, sm)
40   output repconnsr(r').remove({m | m = ⟨STATE, ...⟩...})
41   output repconnsr(r').send(sm)

42 upon global call isCorrectStateMessage(sm, r' ∈ ℛ, state) do
43   return sm = ⟨STATE, r', state⟩τ(r', ℳ, 0, 0) ∧ hasValidCertificate(sm)
      ∧ isCorrectState(state)

44 upon input call receive(sm = ⟨STATE, r', state⟩...) do
45   if isCorrectStateMessage(sm, r', state) then
46     output vmr.installState(state)

```

References

- [1] D. J. Armstrong. The quarks of object-oriented development. *Communications of the ACM*, 49(2):123–128, 2006.
- [2] J. Behl, T. Distler, and R. Kapitza. Hybster — A highly parallelizable protocol for hybrid fault-tolerant service replication. <http://publikationsserver.tu-braunschweig.de/get/64440>.
- [3] J. Behl, T. Distler, and R. Kapitza. Hybrids on steroids: SGX-based high performance BFT. In *Proceedings of the 12th European Conference on Computer Systems (EuroSys '17)*, 2017.
- [4] C. Cachin, R. Guerraoui, and L. Rodrigues. *Introduction to Reliable and Secure Distributed Programming*. Springer Publishing Company, Incorporated, 2nd edition, 2011.
- [5] M. Castro. *Practical Byzantine Fault-Tolerance*. PhD thesis, MIT, 2001.
- [6] M. Castro and B. Liskov. A correctness proof for a practical Byzantine-fault-tolerant replication algorithm. Technical report, Cambridge, MA, USA, 1999.
- [7] M. Castro and B. Liskov. Practical Byzantine fault tolerance. In *Proceedings of the 3rd Symposium on Operating Systems Design and Implementation (OSDI '99)*, pages 173–186, 1999.
- [8] M. Castro and B. Liskov. Practical Byzantine fault tolerance and proactive recovery. *ACM Transactions on Computer Systems*, 20(4):398–461, Nov. 2002.
- [9] C. Dwork, N. Lynch, and L. Stockmeyer. Consensus in the presence of partial synchrony. *Journal of the ACM*, 35(2):288–323, Apr. 1988.
- [10] M. P. Herlihy and J. M. Wing. Linearizability: A correctness condition for concurrent objects. *ACM Transactions on Programming Languages and Systems*, 12(3):463–492, July 1990.
- [11] N. A. Lynch. *Distributed Algorithms*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1996.
- [12] F. J. Meyer and D. K. Pradhan. Consensus with dual failure modes. In *Seventeenth International Symposium on Fault Tolerant Computing*, volume 2, pages 214–222, Apr. 1987.
- [13] F. B. Schneider. Implementing fault-tolerant services using the state machine approach: A tutorial. *ACM Computing Surveys*, 22(4):299–319, 1990.
- [14] P. Thambidurai and Y. keun Park. Interactive consistency with multiple failure modes. In *Proceedings of the 7th IEEE International Symposium on Reliable Distributed Systems (SRDS '88)*, pages 93–100, Oct. 1988.